Screening Through Soft Spending Limits: Evidence from the Medicare Therapy Cap^{*}

Ashvin Gandhi

UCLA & NBER

Maggie Shi University of Chicago & NBER

April 18, 2025

Abstract

Governments and firms often employ soft spending limits to restrict overspending while still allowing exceptions on a case-by-case basis. This paper studies a Medicare policy which capped per-patient physical therapy spending, with exceptions for patients with documented medical need. The cap reduced spending by 8 percent without harming patient health, with the targeting improvements driven by Medicare discretion in granting exceptions rather than improved provider screening. However, the documentation requirement also introduced horizontal inequity: conditional on need, lower-income and minority patients were more likely to be screened out, as they tended to see providers with poorer documentation practices.

^{*}We thank Zhijian Li for excellent research assistance, as well as Jason Falvey, DPT; Brian McGarry, PhD, PT; Patti Marquardt, PT; and Xiao Zheng, MD for sharing their clinical expertise. We thank Zarek Brot, Josh Gottlieb, Tim Layton, Riley League, Jetson Leder-Luis, Ryan McDevitt, Dan Sacks, Jacob Wallace, and seminar participants at ASHEcon, the Junior Health Economics Workshop, the Toulouse School of Economics, IU Bloomington, Chicago Health Economics Workshop, the Berkeley Health Economics Workshop, UIUC, the Hoover Institute, UC Santa Barbara, Stanford, BFI Health Conference, NBER SI Economics of Health meeting, Emory, the Midwest Health Economics Conference, Monash University, the American Health Econometrics Workshop, the Southeastern Health Economics Study Group, the Penn Leonard Davis Institute, UC Chicago Harris, and Bowdoin College for helpful comments. The authors gratefully acknowledge support from the National Institute on Aging (#T32-AG000186) and Arnold Ventures. All errors are our own.

1 Introduction

Governments and firms frequently delegate spending decisions to agents, but these agents may choose to spend in ways the principal deems wasteful. One way to mitigate this is to require agents to seek approval for their spending decisions by submitting documentation that justifies their expenditure. When these requirements apply above a stated threshold, they form a spending limit that is "soft" in that can be bypassed on a case-by-case basis. For example, government procurement policies often require documentation in order to exceed statutory cost caps (CFR, 2020). More broadly, requiring justification for approval to receive *any* funds—effectively placing a limit at zero—is a common feature of many expenditure programs, including those providing disability benefits (SSA, 2008), unemployment insurance (DOL, 2025), and disaster assistance (FEMA, 2021).

These policies aim to curb wasteful spending while still allowing flexibility for exceptional needs above the limit. They can improve efficiency by screening in two ways: the documentation can inform the principal's decision to *deny* or approve a request, and the associated ordeal can lead agents to "self-screen" by *deterring* low-value requests (Nichols and Zeckhauser, 1982). However, they may also inadvertently favor agents with greater administrative capability: those who are better at producing documentation will be less deterred and face fewer denials. This can generate horizontal inequity since agents' administrative capacities do not necessarily reflect the value of their requests. Despite the widespread use of such "soft" screening mechanisms, there has been little empirical work unpacking whether and how they improve targeting, as well as the potential equity-efficiency tradeoff they impose.

This paper studies these questions in the context of a spending limit imposed by insurers on healthcare spending. While initial care decisions are delegated to patients and their providers, the insurer ultimately has the final say in approving or denying reimbursement for care. One way insurers mitigate overspending is by setting per-patient spending limits on certain services and requiring providers to submit documentation to justify going above the limit.¹ This practice is especially common for treatments susceptible to moral hazard, such as for physical therapy, psychotherapy, chiropractic, and dental services.²

¹Examples include visit notes, test results, imaging reports, and evidence of prior treatments.

 $^{^{2}}$ In the extreme, insurers can place restrictions on the *first* unit of care, a practice known as prior

Specifically, we study a spending limit imposed on physical therapy (PT) within the Medicare program, where policymakers were concerned about medically unnecessary overuse of services like manual massage and supervised strength training. The "therapy cap" was an annual, per-patient spending limit of \$1,740, or about 11 weeks of care. Providers could request exceptions for individual patients to exceed the limit based on documented medical need. Medicare then approved or denied these requests based on whether the available information indicated further care was medically necessary.

Our setting of a soft spending limit in healthcare is well-suited for understanding how soft screening mechanisms work more generally. The policy we study imposes a sharp threshold on a recurring service, which allows us to precisely identify not only *who* is screened out, but also *why*. We can track patients as they approach the limit and focus on the ones who stop just below it. Among these patients, we can distinguish between the two mechanisms driving savings: deterrence—patients who don't attempt to go past the limit—and denials patients who attempt to exceed it, but are stopped at the insurer's discretion. Furthermore, observing health outcomes allows us to evaluate the extent to which the cap targets medically unnecessary care. Finally, our data include measures of documentation use, allowing us to assess the role of provider administrative capacity in mediating these effects.

We first quantify the savings from the therapy cap with a difference-in-bunching estimator that compares the distribution of PT spending before and after the 2006 introduction of the cap. Overall, the therapy cap reduced spending by eight percent relative to pre-reform levels.³ Fifty-eight percent of savings are due to Medicare denying requests to exceed the cap, while the remaining 42 percent stems from providers being deterred from making attempts.

We then employ two approaches to evaluate how well these spending reductions targeted medically unnecessary care. The first is a patient-level analysis that leverages the fact that the cap resets each calendar year and is therefore more binding for patients who start PT earlier in the year. We use this variation to isolate spending reductions that are directly attributable to the cap, and then compare outcomes for patients more or less affected by the cap. Our results indicate that the cost savings from the cap were well-targeted, as they

authorization (Dillender, 2018; Brot-Goldberg et al., 2023).

³In contrast, a 1999 "hard" cap without exceptions reduced spending by 14 percent.

neither resulted in substitution to alternatives like opioids, pain procedures, and orthopedic surgeries, nor in worsening health outcomes like hospitalizations, emergency department visits, and nursing home stays.

Our second approach then looks at patients as they approach the cap or attempt to go over it, and tests whether the patients screened out by the cap have comparatively low medical need, as indicated by their prior utilization and diagnoses. Importantly, we separately consider who attempts conditional on approaching the cap, and who is approved conditional on making an attempt. This allows us to distinguish between screening done by Medicare through the denial of attempts versus by patients and providers through the deterrence of attempts. We find that while the savings can be attributed to a combination of deterrence and denials, the improvements in targeting on need appear to be driven entirely by Medicare's denials. The fact that denials align with medical need means that Medicare's enforcement behavior was consistent with its stated objective of reducing medically unnecessary care. However, the lack of similar targeting through deterrence runs counter to the prediction that providers should respond to the associated ordeal by "self-screening" more (Zeckhauser, 2021). Instead, it suggests that the ordeal acts as a fairly blunt screening tool, reducing spending in an untargeted way.

Beyond medical need, the cap may also have inadvertently screened along other dimensions. Given that the cap introduced a documentation requirement that providers had to meet, we test whether it also screened on providers' administrative capacity to comply with this requirement. We use provider size as a proxy for administrative capacity (League, 2022; Dunn et al., 2024), and find substantial disparities by provider size in Medicare denial rates. Conditional on requesting an exception to the cap, a patient with median need faces an 80 percent (12.7 percentage point) lower denial rate when receiving care from a large provider rather than a small provider. Despite this, smaller providers are not more deterred, suggesting that they are making mistakes in submitting requests with low probability of approval or that their unsuccessful attempts are not as costly.

Because in our setting smaller providers tend to see minority and low-income patients, these differences in denials across providers translate into meaningful horizontal inequity by patient race and income. Non-white patients, Medicaid enrollees, and Part D Low Income Subsidy recipients with median need respectively experience denial rates that are 43 percent, 23 percent, and 21 percent higher than their counterparts. Furthermore, we show that these disparities are driven by patient sorting across providers, as they largely disappear with the inclusion of provider fixed effects.

Finally, we investigate the mechanisms driving these differences in denial rates across large and small providers. We show that whether an exception request is approved or denied is closely tied to whether the provider has supporting documentation available, and that larger providers are much more likely to comply with documentation requirements. We decompose potential sources of this size advantage in compliance and find that much of it comes from learning-by-doing: providers are more likely to use documentation as they gain experience navigating the cap, and large providers accrue experience more quickly simply by having more patients.

This paper contributes to the literature on the effectiveness and targeting of screening mechanisms (Nichols and Zeckhauser, 1982; Kleven and Kopczuk, 2011; Alatas et al., 2016; Deshpande and Li, 2019; Finkelstein and Notowidigdo, 2019; Lieber and Lockwood, 2019; Ida et al., 2022; Brot-Goldberg et al., 2023; Shepard and Wagner, 2024). Many of these mechanisms screen in two stages: first via a deterrence stage where agents screen themselves out, and then through a denial stage where the regulator makes the final decision. Our paper is the first to distinguish between the two channels and evaluate their savings and targeting properties separately. We find that making this distinction matters: while deterrence and denials both contribute to savings, they have very different targeting effects. The beneficial targeting comes entirely from the regulator's denial decisions, rather than agents self-screening in response to the ordeal. This suggests that reducing the administrative burden imposed on the agent by documentation requirements while still preserving their informational content — for example, through automated electronic reporting — could be welfare-improving.

We also speak to the literature on disparities in takeup of public programs (Currie, 2006). Much of this literature has focused on barriers to beneficiaries like information gaps, stigma, and application costs. We highlight a largely unexamined factor: the intermediaries who interact with public programs on behalf of beneficiaries and, more specifically, these intermediaries' administrative ability. In our setting, providers serve as key liaisons between patients and Medicare. Beyond just providing care, they also handle the administrative paperwork that can shape a patient's access to the Medicare program. In this way, their role is analogous to that of tax preparers (Kopczuk and Pop-Eleches, 2007; Zwick, 2021), disability lawyers (Hoynes et al., 2022), social work case managers (Evans et al., 2024), and mortgage brokers (Woodward and Hall, 2012). We find that heterogeneity in administrative ability across these intermediaries—attributable mostly to differences in their accumulated regulatory experience—leads to substantial variation in whether their clients successfully qualify for benefits. When combined with uneven sorting across intermediaries, this translates into meaningful horizontal inequity across beneficiary race and income, even among beneficiaries with identical observable need. As these disparities arise from differences in intermediaries, this suggests that interventions to increase provider awareness of new regulation, like training programs or administrative support, are likely to be more effective at reducing horizontal inequity than efforts focused on relieving beneficiary-facing barriers.

Finally, we contribute to the literature on policies aimed at encouraging providers to reduce wasteful healthcare spending. This includes tools like prior authorization (Dillender, 2018; Brot-Goldberg et al., 2023; Eliason et al., 2024), audits or denials (Macambira et al., 2022; League, 2022; Shi, 2024), and anti-fraud enforcement (Howard and McCarthy, 2021; O'Malley et al., 2023; Leder-Luis, 2023). We add to this literature by examining spending limits as a policy tool to curb wasteful care specifically on the intensive margin. Our data and policy setting allow us to closely examine the underlying mechanisms at play, particularly on the distinct effects of deterrence and denials, as well as on the key role provider administrative ability plays in mediating these effects.

The paper proceeds as follows. Section 2 discusses the policy context and places it within a conceptual framework, and also describes the data we use and measures we construct. Section 3 presents estimates of the overall effects on savings and patient health. Section 4 characterizes the targeting properties of the cap as well as horizontal inequity. Section 5 explores the mechanisms driving the horizontal inequity, and Section 6 concludes.

2 Policy Context and Data

2.1 Outpatient Physical Therapy Care in Medicare

In 2017, Medicare Part B spent \$4.9 billion on outpatient physical therapy services for 2.6 million beneficiaries, or about 7 percent of all traditional Medicare beneficiaries. Medicare beneficiaries can receive three main types of therapy in the outpatient setting: physical therapy (PT), occupational therapy (OT), and speech-language pathology (SLP).⁴ PT accounts for over three-quarters of total outpatient therapy spending, and includes the diagnosis, management, and prevention of physical dysfunction and pain. Within the Medicare population, the most common procedures associated with PT are therapeutic exercises, manual therapy, and electronic stimulation (Table A1). The most common associated diagnoses are for low back pain, shoulder pain, and lower leg pain.

A patient is typically referred to PT by a physician, and the physical therapist must develop a plan of care for the patient, which the physician certifies. ⁵ The plan of care includes a diagnosis, long-term treatment goals, and the type, quantity, duration, and frequency of therapy services. Outpatient Medicare PT services are primarily provided in private practices (34%), nursing facilities (38%), and hospitals (15%) (APTA, 2020). Regardless of setting, outpatient PT is subject to the standard Medicare Part B cost-sharing rules: for approved services, Medicare pays 80 percent of allowed charges and the patient is responsible for the remaining 20 percent.

Patients might like to continue treatments such as directed exercise and massage even after their clinical benefits do not justify the costs for Medicare (OIG, 2016, 2017). As such, policymakers have long been concerned about medically unnecessary overuse of therapy services in the Medicare program. In 1993, the predecessor to the Centers for Medicare and Medicaid Services (CMS) launched a task force to address widespread reports of overbilling

⁴Occupational therapy can resemble physical therapy—e.g., OT frequently also frequently entails therapeutic exercises and manual therapy — but is differentiated in its focus on improving a patient's ability to perform activities of daily living like bathing and dressing. SLP aims to improve patient's ability to speak and swallow.

⁵It is also possible for patients to be evaluated by a physical therapist before seeing a physician ; however the physical therapist still needs to obtain a physician referral or signature at some point in order to receive payment for these evaluation services.

for therapy (US GAO, 1996). In the years that followed, CMS tried many policies to curb therapy spending, including the spending limits we study. An Office of the Inspector General (OIG) audit found that a third of Medicare outpatient therapy claims were still for medically unnecessary services, and over half of claims reviewed were not compliant with medical necessity, coding, or documentation requirements (OIG, 2018). According to the OIG, the most common reason for unnecessary PT services was patients receiving excessive amounts of therapy—that is, although a patient may have initially needed therapy, the amount, frequency, or duration they received went beyond "standards of practice."⁶ Thus, efforts to curb excessive PT spending have focused on limiting the amount of per-patient spending in the form of annual "therapy caps."

2.2 Medicare Therapy Cap

Medicare's policies on physical therapy spending can be divided into three regimes: a oneyear "hard cap" in 1999, a six-year period with effectively no cap, and a period with a "soft cap" beginning in 2006 that continues to today. We focus primarily on the 2006 soft cap, but discuss the 1999 hard cap in Appendix Section C. The soft cap was first implemented as two \$1740 caps in January 2006 — one placed on PT and SLP, and another placed on OT.⁷ Medicare introduced a process through which providers could request exceptions for medically necessary services above the cap, thus making it "soft." When billing for services that would push a patient above the cap, they were supposed to indicate that they had documentation justifying medical necessity by using a billing code called the "KX" modifier code (CMS, 2006a). Providers did not have to attach the documentation to their claim; instead, using the modifier indicated that they "attested" that the documentation indicated that spending above the cap was "reasonable and necessary," and was available should Medicare request it. All attempts to bill over the cap were supposed to include this modifier

⁶Medicare requires the following for a therapy service to be considered "reasonable and necessary": (1) that the services are an effective treatment for a patient's specific condition, (2) that the service must be performed (or supervised) by a therapist, (3) that the service is expected to improve a patient's condition or is necessary to maintain their condition, and (4) that the amount, frequency, and duration of therapy follow standards of practice (CMS, 2020).

⁷We limit our analysis to the PT/SLP cap given that PT accounts for the majority of outpatient therapy spending.

code, though enforcement of this rule was inconsistent—16 percent of claims in which the modifier documentation should have been used but was not were still paid out in full in 2006.

In using the documentation modifier code, the provider attested that the following documentation was available for review by Medicare: evaluation and plan of care, certification and re-certifications with evidence of physician (or non-physician practitioner) approval, progress reports, treatment notes, as well as potential separate justifications to indicating the reasoning for services that are "more extensive than is typical for the condition treated" (CMS, 2024a). The documentation is expected to justify that the patient requires the skill of a therapist, that the services are the appropriate type, frequency, intensity, and duration for the particular needs of the patient, as well as provide relevant information about the patient's functional abilities (CMS, 2006a).

The cap was enforced by Medicare through claim denials, as shown by the increase in denial rates at the cap that appears in 2006 in Figure B1. In order for the therapist to charge the patient for care in the event of a Medicare denial, she must issue an "advance beneficiary notice" (ABN) *prior* to the delivery of the service. The ABN informs the patient why Medicare may not pay for a specific claim and allows them to choose whether to go forward with the service and, in the event of a denial, to accept financial responsibility. If a claim is denied without the issuance of an ABN, then the therapist is not allowed to charge the beneficiary and is financially liable for the cost of any services rendered (CMS, 2013, 2018). Medicare is required to issue an initial determination of denial within 30 days of receipt of the claim, though anecdotally they often reach this decision much faster (CFR, 2009).

Figure 1 summarizes the actions and outcomes at the therapy cap for patients, providers, and Medicare. In evaluating the cap's targeting properties, we consider two key decision points in a patient's PT care trajectory: the week in which they "attempt" to bypass the cap via an exception, and the week in which they "approach" the cap ahead of a potential attempt. As a patient approaches the cap, they make a decision with their provider of whether to continue care or not. Patients who approach the cap but stop without an attempt are considered "deterred" by the cap. We interpret this deterrence as a response to the ordeal associated with the cap—this includes the costs of care, cost of documentation, and the uncertain reimbursement net of denials. Those who are not deterred then continue care, produce documentation, and attempt to bill above the cap. Medicare then decides whether to approve or deny these requests. Patients who attempt to exceed the cap but are stopped by Medicare are classified as "denied," while those who continue past the cap are classified as "approved." Section 2.3 describes in further detail how each of these are identified and defined in our data.



Figure 1: Diagram of Therapy Cap Actions and Outcomes

Footnote: This figure illustrates patient, provider, and Medicare's actions and outcomes at the therapy cap as discussed in Section 2.2.

In sum, patients who bunch under the cap can be categorized as either those who approached and were deterred or those who attempted and were denied. Making this distinction between deterrence and denial is important because they capture two different channels through which the cap operates. Through the deterrence channel, providers and patients decide whether to try to continue care after assessing their private costs and benefits. Through denials, Medicare plays an active role in deciding who is allowed past the cap. We empirically estimate the size of each channel and characterize their targeting properties separately.

2.3 Data and Sample

Data Our main source of data is the Medicare 20% Carrier claims files. These data include all Part B office-based spending for a random 20 percent subset of traditional fee-for-service Medicare beneficiaries. At the line-level, the key variables are procedure (HCPCS) codes, units, and final payments. Lines can be aggregated into claims, which each correspond to a single PT visit. At the claim-level, the key variables are diagnosis codes (ICD-9), billing modifier codes (including the KX modifier), dates of service, provider identifiers, and an indicator for payment denial. We define a "provider" as a unique combination of tax identification number (TIN) and state.⁸ We use TINs to define providers as opposed to National Provider Identifiers (NPI) for two reasons. First, providers were only required to report NPI in 2008, which is after our analysis period.⁹ Second, we expect that many of the behaviors or investments providers would take up in response to the cap would be implemented at a *practice*-level, as opposed to at the individual provider level. These include, for example, any upgrades electronic medical record systems, improvements to documentation standards, or changes in screening practices.

We supplement the Carrier claims with additional spending and utilization measures from the 20% Medicare Provider Analysis and Review (MEDPAR) and 20% Outpatient files. We also use information on patient demographics, chronic conditions, prior utilization, and mortality from the Medicare Master Beneficiary Summary Files (MBSF). Finally, we use 2006 zip-level income statistics reported by the Internal Revenue Service Statistics of Income (IRS, 2025).

Sample Definition We then narrow down to a sample of patients who receive in-office physical therapy.¹⁰ We follow Amico et al. (2015) in identifying PT claims via the HCPCS code and the PT modifier code. We focus on the 91 percent of these patients who only see one provider for PT the entire year. Among these patients, we limit to patients with "regular" PT: those that have at least 5 weeks with at least \$50 spending a week in the

⁸We count TINs that operate in separate states as separate providers in order to capture the notion of a physical practice. 98% of TINs operate only in one state.

 $^{^{9}}$ In 2008, the average TIN is our sample is associated with 2.5 NPIs.

 $^{^{10}}$ The cap technically applied to PT *and* speech therapy combined, but we remove the 4 percent of patients who ever receive speech therapy.

calendar year. We do this to eliminate outliers of patients who have short but expensive PT episodes (e.g., one week of over \$1000 in spending with no spending in other weeks) due to concerns that these reflect misreporting or that these patients are not comparable to the rest of the sample.

We then create four different samples for each of our analyses; we summarize the sample definitions broadly here and describe them in further detail in the respective sections. The bunching analysis in Section 3.1 restricts to patients with end of year spending within [-\$800,\$1600] of the cap, the health analysis in Section 3.2 restricts to patients with an injury diagnosis and no PT in the six months prior to their first session, and the screening analysis in Section 4 restricts to patients who "approach" and "attempt" the cap (as defined below). Table A2 reports summary statistics for each sample in 2005 and 2006. Columns 1 and 2 show the summary statistics for the bunching analysis sample, columns 3 and 4 report statistics for the health analysis sample, and columns 5-8 report statistics for the "approach" and "attempt" samples. Overall, within each sample the 2005 and 2006 subsamples are relatively balanced on measures of demographics and prior utilization. While the total number of visits and amount of PT spending typically decreases from 2005 to 2006, the samples are fairly balanced in terms of the average per-visit PT spending. The number of patients in the bunching and health samples increases from 2005 to 2006, which follows a pre-reform trend of increasing numbers of patients receiving PT in this time period. Looking closer to the cap, there is a slight drop in the number of patients who approach the cap and a notable decrease in the number who make an attempt from 2005 to 2006.

Denials and Attempts Denials can be defined at the claim-, patient-, and attempt-weeklevel. We classify a *claim* as "denied" if the Carrier claim payment denial code indicates the claim was denied.¹¹ Ninety-four percent of denied claims for PT in 2006 are associated with no payment, while many of the remaining claims largely consist of partial payments.

Given the prevalence of denials in our sample, we distinguish between two ways to summarize the spending associated with a claim: the *billed* and the *paid* amount. The amount paid reflects the payment the provider receives after the bill has been processed, net of de-

¹¹See ResDac (last accessed 3/8/24) for more detail on this code.

nials. This is the final payment amount reported directly in the claims, and is also what the therapy cap applies to. The billed amount reflects the care that was actually provided and that the provider demanded payment for, prior to any denials. Since we do not observe the amount of payment demanded for denied line items, we construct it using the procedure code. Specifically, we impute the billed amount using that provider-patient pair's average per-unit payment on claims without a denial for that procedure code.¹²

We then use this distinction between billed and cumulative amounts to identify the weeks in which a patient attempts to go over the cap, following Figure 1. Weeks with "attempts" are defined as weeks in which the cumulative paid amount up to the prior week was below the cap, and the billed amount in that week plus the prior week's cumulative paid amount is projected to go over the cap. We classify an *attempt week* as "denied" if at least one of its associated claims has a denial. Conversely, an attempt week is considered "approved" if none of the associated claims have a denial and the patient continues past the cap.

Finally, we classify a *patient* as "denied" if they make an attempt but never successfully make it over the cap that year—that is, their cumulative paid amount at the end of the year is below the cap. Conversely, a patient is "approved" if they end the year with a cumulative paid amount above the cap.

Weeks from Cap We also assign to each week how many "weeks from the cap" the patient is. We only consider weeks in which a patient received PT care.¹³ When "weeks from cap" measure is negative, it denotes how many additional weeks of care a patient would have to receive to reach (and pass) the cap. When "weeks from cap" is (weakly) positive, it denotes how many weeks ago the patient passed the cap.

The method for assigning "weeks from cap" differs depending on whether a patient's cumulative paid spending at the end of the year fell above or below the cap, and if they are ever observed making an attempt to go over the cap. For patients who end the year above

 $^{^{12}}$ If that provider-patient pair has never successfully billed for that procedure code, we use the provider's average payment. If the provider has never successfully billed for that procedure code, we use the average payment for that code among providers with the same Medicare Administrative Contractor. This approach is similar to that of Dunn et al. (2024).

¹³The week-to-week visit patterns in the data suggest that care is scheduled on a weekly basis. Figure B2 shows that PT care tends to be scheduled on the same day every week, as well as on a Monday/Wednesday/Friday or Tuesday/Thursday cadence.

the cap, week 0 is defined as the week of their attempt, and all other weeks with positive PT spending are defined relative to that week.¹⁴ For patients who make an attempt but never get successfully get past the cap (thus ending the year below the cap), their last week is considered week -1, and all other weeks are defined relative to that. Finally, for patients who never make an attempt and end the year below the cap, we extrapolate based on their prior weekly spending to calculate how many "weeks away" the patient is from reaching the cap. Specifically, we calculate the patient's 5-week rolling average of PT spending, take the difference between the cap value and the cumulative amount billed that week, and divide by the maximum of the patient's 5-week average or the average weekly spending in the sample.¹⁵ Figure B1 illustrates week-level denial rates by "weeks from cap," and shows a clear increase at week -1 in denials that begins in 2006.

Approaches After defining weeks from cap, we can then define the week in which the patient "approaches" simply as the week where the patient is -1 weeks away from the cap. For patients who make an attempt, this means that their approach week is just the week before their first attempt. For patients who never make an attempt, it is the point at which one more week of their usual care is extrapolated to take them over the cap.

Documentation Indicator While we cannot directly observe documentation in the claims data, we proxy for it by looking for whether the "KX modifier" code is present on a claim. As discussed in Section 2.2, this modifier code was introduced as part of the 2006 therapy cap and is used by providers to indicate the availability of documentation justifying the medical necessity of spending over the therapy cap. Table 1 shows that having documentation, as indicated by the presence of this modifier code on at least one claim in an attempt week, substantially increases the likelihood that that week's attempt is approved. This correlation is robust to the inclusion of controls for patient and provider size, the R^2 increases 2-4-fold

 $^{^{14}}$ If a patient ends up past the cap but only after making multiple attempts, we consider their *first* attempt to be week 0.

¹⁵We divide by the maximum of the two to ensure that patients with low prior weekly spending are never implausibly far away from the cap (e.g., over 52 weeks away from the cap). Dividing by the maximum of the two effectively left-censors the "weeks from cap" measure at -8.

once this indicator is included in the regression.

Predicted PT Spending We also use the claims data to construct a measure of ex ante patient clinical need for PT. Using data from pre-reform years, we predict what a patient's 12-month PT spending would be absent the therapy cap, given their spending and utilization patterns prior to their first PT. We implement this using gradient-boosted decision trees from the LightGBM package. The predictors are age, sex, utilization and spending in the previous calendar year available in the MBSF Cost and Utilization file (in-office spending, Part B drug, outpatient procedure, inpatient, testing, imaging, hospice, evaluation and management, durable medical equipment, dialysis, and other), chronic conditions at the end of the previous calendar year, PT and OT spending in the previous calendar year, inpatient and SNF stays within the last 6 months (spending, number of visits, Diagnosis Related Group of their most recent inpatient stay, length of stay of their most recent inpatient stay, and days since last visit), in-office spending in the last six months, and an indicator for having an auto exception diagnosis in the last 6 months.¹⁶ The model is trained on patients who approach or attempt the cap in 2004 and 2005, prior to the implementation of the therapy cap. Thus, we are predicting what a patient *would* have spent on PT, absent the cap. We then apply this prediction to patients who approach and attempt the cap in 2005 and 2006. We discuss the machine learning methodology and model fit in further detail in Appendix Section D.

3 Overall Effects of the Cap

3.1 Savings

Methodology and Sample Construction To quantify the Medicare savings from the cap, we apply methods from the "bunching" literature (Kleven, 2016). In our context, the pre-policy distribution serves as a natural counterfactual. In particular, we will compare distribution of Medicare PT spending in 2006 to the pre-reform distribution in 2005, restricting

 $^{^{16}}$ Importantly, note that the model is *not* trained on factors we later consider in our test for horizontal inequity: race and income (as well as zip code, which could be a proxy for both).

to the area around the 2006 cap—the "manipulation region." Medicare pays for 80 percent of allowed charges and patients are responsible for the remaining 20 percent, so the cap appears at $0.8 \times \$1740 = \1392 in the distribution of per-patient Medicare PT spending. We restrict to a region from \$800 below to \$1600 over this amount; this region has been chosen such that the "missing mass" to the left of the cap is approximately equal to the "excess mass" to the right. We interpret the difference in spending between the two distributions as the savings implied by the cap.

Rather than directly comparing the 2005 and 2006 distributions, we adjust the 2005 spending calculation to account for changes in procedure prices and in the total number of patients receiving PT from 2005 to 2006. Let $\bar{r}(F_j, p_k)$ be the per-patient spending in the manipulation region using the distribution in year j and prices in year k:

$$\bar{r}(F_j, p_k) := \int \left(q \cdot p_k\right) dF_j(q),$$

where F_j is a distribution in the region around the cap in year j over quantity q for each procedure, p_k is a vector of prices for each procedure in year k, and the integral is taken over all PT procedure codes. Thus, $\bar{r}(F_{06}, p_{06})$ denotes the actual average spending around the cap in 2006, and $\bar{r}(F_{05}, p_{06})$ denotes the average spending under the 2005 distribution, price-adjusted for 2006.¹⁷

In order to convert the difference between these two averages into a measure of total savings, we multiply the per-patient spending by the number of patients in this region in 2006, denoted as N_{06} . This ensures that the savings calculation is not affected by secular changes in the total number of patients receiving PT in this period. Thus, the savings in 2006 dollars is defined as:

$$S_{06} := N_{06} \left(\bar{r}(F_{05}, p_{06}) - \bar{r}(F_{06}, p_{06}) \right),$$

¹⁷To impute prices, we use the *provider*-level 2006 price to account for geographic adjustments or any other provider-specific idiosyncrasies that affect the price they receive per procedure. For procedures that a provider rendered in 2005 but not 2006, we replace the 2005 price with the average 2006 price for that procedure in the same Medicare Administrative Contractor (MAC) region to account for any geographic differences in Medicare prices (League, 2022).

and the savings as a percent of 2005 spending is defined as:

$$S_{06}^{\%} := \frac{\bar{r}(F_{05}, p_{06}) - \bar{r}(F_{06}, p_{06})}{\bar{r}(F_{05}, p_{06})}.$$

Results Figure 2 depicts the 2005 and 2006 spending distributions and their difference in the range from \$800 below the cap to \$1600 over the cap. The sum of the difference between the two distributions is \$83 million, or 7.6 percent of 2005 spending in the region around the cap. We estimate that there were 532,000 additional denials in this region in 2006, which implies that the return to Medicare per additional denial due to the cap was \$156.

Figure 2: Distributions of Spending Around Cap in 2005 and 2006



This figure plots (a) the distributions of end-of-year physical therapy spending around the cap in 2005 and 2006 and (b) the difference in the distributions between 2005 to 2006. Distance from cap is calculated in bins of \$50 relative to the 2006 cap and shares are calculated as the share of patients within [-\$800, \$1600] of the cap. Data: 20% Medicare Carrier claims.

Interpreting the difference between the 2006 and 2005 distributions as the savings from the cap requires two key assumptions. The first is that the 2005 distribution captures what the 2006 distribution would have been absent the reform. This would be violated if this part of the distribution is not stable from year to year, in the absence of changes in Medicare policy. We can validate this assumption by comparing the 2005 distribution to the 2004 distribution, neither of which was subject to the therapy cap. Figure B3 shows that when comparing two consecutive years without a therapy cap, there is little difference in the spending distributions in this region.

The second assumption is that the introduction of the therapy cap did not change the share of patients who ended up inside or outside of the manipulation region around the cap. In other words, the only reason $N_{06} \neq N_{05}$ is due to secular trends over time in the total number of patients receiving PT and not because of the therapy cap. This allows us to normalize patient count across the two years by multiplying by the 2006 patient count in this region. For this to be true, we need that there was no "extensive margin" response to the cap which differentially drew patients into or out of the manipulation region. It would be violated if, for example, providers who cut back on care for patients above the cap now have more capacity and use this additional capacity to accept more patients who end up far below the cap. In that case, the relative share of patients outside of the manipulation region would increase as a response to the cap.

We evaluate this second assumption in two ways. First, looking at the full distribution of spending in 2005 and 2006 in Figure B4 panel (a), we note that the lower part of the distribution outside of the manipulation region is relatively stable. There does not appear to be any marked changes in the share of patients below the manipulation region. Second, if there were an extensive margin response, we would expect that it would be driven by providers who had relatively more patients over the cap in the pre-period, as they would be the ones experiencing the largest capacity expansions. Thus, increases in lower-spend patients should be concentrated among these providers. We explore this in Figure B4 panel (b), which plots average patient count over time, splitting by providers who had high vs. low shares of patients over the cap in 2005. Providers with high 2005 over-cap shares do not appear to see more patients starting in 2006, indicating that they did not respond to the cap by taking on more patients at the lower end of the spending distribution.

Comparison to a Hard Cap To give the savings from the 2006 soft cap more context, we also compare it to the hard cap which was implemented in 1999. The difference between the two regimes is that there was no exceptions process in 1999 for patients to get care above the cap—effectively, all attempts to go past the cap were denied. The policy context for the 1999 cap is discussed in greater detail in Appendix Section C. Due to data limitations, we

cannot compare 1999 to a pre-period year but instead compare it to 2000, after the cap was repealed. If there are any lingering effects of the 1999 cap in 2000, then this would bias our savings estimate downward.

Figure B5 plots the 1999 and 2000 spending distributions and differences in distributions in the range from \$700 (in 1999 dollars) below the cap to \$1300 above the cap, which is equal to approximately \$800 and \$1600 in 2006 dollars. Here, the savings are defined as:

$$S_{99} := N_{99} \left(\bar{r}(F_{00}, p_{99}) - \bar{r}(F_{99}, p_{99}) \right)$$

Taking the difference between the 1999 distribution and the price-adjusted 2000 counterfactual, we find that the hard cap reduced spending by 14.1 percent (\$23.4 million in 2006 dollars). Compared to the 7.6 percent reduction from the 2006 cap, this confirms that the 2006 therapy cap was indeed "soft"—the savings are diminished by almost half when providers are allowed to request exceptions to exceed the cap.

Quantifying the Deterrence and Denial Channels Figure 2 shows that the therapy causes patient spending to bunch at the cap. To unpack what is driving this bunching, in Figure 3 we plot the distribution of end-of-year PT spending by the number of *weeks* a patient is from the cap. Panel (a) plots the overall difference between the 2005 and 2006 distributions by weeks to the cap. The excess mass to the left of the cap is concentrated in the last 3 weeks leading up to the cap, and there is a 93 percent (8.1 percentage point) increase in the share of patients who stop one week before the cap. This excess mass to the left of the cap, with effects extending even up to 10 weeks away from the cap.



Figure 3: Distribution in Weeks to Cap, 2005-2006

(a) Difference between 2005 and 2006 (b) D distributions in weeks to cap

(b) Difference in weeks to cap, split by reason for stopping at cap

This figure plots the difference in the 2005 and 2006 distributions of cumulative physical therapy. Distance from the cap is calculated in terms of weeks of care relative to the 2006 cap and shares are calculated as the share of patients within [-5,10] weeks of the cap. Panel (a) shows the overall difference between 2005 and 2006, and panel (b) splits the difference between patients with an attempt to go past the cap in their last week but were denied ("Attempt and denied", red), patients with an attempt who were not denied ("Attempt and successful," orange), and patients with no attempt ("Never attempt", green). Section 2.3 describes the construction of the weeks from cap measure. Data: 20% Medicare Carrier claims.

Figure 3 panel (b) then decomposes the excess and missing mass into three categories. As in panel (a), we show the difference between the 2005 and 2006 distributions. Patients with an attempt to go over the cap can either be classified as "Attempt and Denied," meaning they made an attempt to bill past the cap but were denied by Medicare, or they are classified as "Attempt and Approved," meaning they ended the year above the cap. The remaining patients are classified as "Never Attempt," and end the year one or more weeks below the cap. Among patients to the left of the cap, we consider them to be "deterred" if they stop without an attempt, and "denied" if they attempt and are denied. Of patients who bunch in the week *immediately* before the cap, 42 percent stop due to deterrence while 58 percent stop because they are denied.

Comparing the distribution of beneficiaries allows us to decompose the share of beneficiaries who stop at the cap into the denial and deterrence channels. However, the decomposition of *savings* attributable to each channel could differ, as this depends not only on the number of beneficiaries who stop, but also their counterfactual spending in the absence of the cap. This counterfactual spending is unobservable, but we can use the distribution of *predicted* spending to estimate the savings decomposition. Specifically, we use the predicted 2004-2005 PT spending measure described in Section 2.3 and Appendix Section D.

Figure B6 plots the distributions of predicted PT spending among "deterred" and "denied" patients who stop one week from the cap in 2005 (panel a) and 2006 (panel b). The savings attributed to each channel can be calculated as the difference between the area under the 2006 and 2005 distributions (panel c). We calculate that the savings from deterrence are \$34 million and the savings from denial are \$51 million, meaning 40 percent of the savings can be attributed to deterrence and 60 percent can be attributed .¹⁸ This split is very close to the share of beneficiaries who are deterred or denied, reflecting the fact that the spreads of the predicted savings distributions are fairly similar between the deterred and denied groups.

3.2 Patient Health Effects

We next consider whether the cap had a direct effect on patient health or led patients to substitute to other forms of care. If there is evidence of such effects, this would imply that the cap-induced savings were poorly targeted. To estimate the causal effect of the therapy cap on health outcomes, we use an instrumental variables (IV) strategy that leverages within-year variation in spending due to the differential "bite" of the cap, depending on when in the year a patient initiates PT. Appendix Section F presents an alternative approach which employs a difference-in-differences estimator and finds similar results.

Methodology and Sample Construction Our identification strategy relies on the fact that the therapy cap resets at the beginning of each calendar year, regardless of when a patient begins receiving PT care. All else equal, the cap is more binding for a patient who starts their care earlier vs. later in a year. As a result, this should generate a negative relationship between when a patient starts PT and their total 12-month PT spending once

¹⁸While most of the excess mass appears exactly one week short of the cap, there is also some excess mass dispersed up to three weeks away from the cap. Once patients who stop up to three weeks away are included in the deterrence channel, then 53 percent can be attributed to deterrence while 47 percent stop due to a denial. Under this categorization, the savings from deterrence rise to \$54 million, or 51% of the total savings.

the cap is in place in 2006. This motivates our instrument: the interaction of the month that a patient starts PT with the treatment year.

We construct our sample by first taking patients who start PT in either 2005 or 2006, where the "starting" PT session is defined by having no PT in the six months prior. In order for the timing of when patients start PT to be a valid instrument for PT spending, we need to assume that the exclusion restriction holds: the only reason for the relationship between when the patient starts PT and their subsequent health outcomes is because patients who start earlier in the year in 2006 are more likely to have their PT spending constrained by the cap.

One potential violation of this assumption would be if the cap changes the composition of who starts care earlier vs. later in the year between 2005 and 2006. For example, if sicker patients were more aware of the cap and strategically timed their care to start later in 2006, this would result in a negative correlation between patient health and PT start date that is not causal. To address this concern, we restrict to a patient population where the PT start date is more likely to be exogenous: those who received an injury diagnosis in the 90 days prior to their first PT visit.¹⁹ We also restrict to patients with 12-month PT spending over \$200 to focus on patients more likely to be affected by the therapy cap.²⁰ As robustness checks, we re-run our analyses on patient populations who may be more vulnerable to spending reductions from the therapy cap—low-income patients and high-need patients. We also look for heterogeneous effects by decile of predicted patient need and by whether a patient had a pain procedure or orthopedic surgery prior to starting PT.

We first estimate the following reduced form specification:

$$Y_i = \sum_{f=1}^{11} \theta_f \ 1(Year_{y(i)} = 2006) \times 1(FirstMonth_{m(i)} = f) + FirstMonth_{m(i)} + Year_{y(i)} + \varepsilon_i,$$
(1)

where θ_f is the coefficient on the interaction of an indicator for 2006 and the first month of a patient's PT, and the omitted month is December. $Year_{y(i)}$ is a year indicator and

 $^{^{19}}$ Appendix Section E describes how injuries are identified in the sample construction.

 $^{^{20}}$ As discussed in Section 3.1 and Figure B4, we find no evidence of an extensive margin response to the cap in the lower part of the spending distribution.

 $FirstMonth_{m(i)}$ is an indicator for the patient's first month. Then the second stage is:

$$Y_i = \beta \widehat{PT_i} + FirstMonth_{m(i)} + Year_{y(i)} + \nu_i, \tag{2}$$

where $\widehat{PT_i}$ is instrumented 12-month PT spending based on Equation 1, scaled so that β can be interpreted as the effect of an additional \$100 in Medicare PT spending on the likelihood of each outcome. Our estimates are clustered by year and start month.

As our health outcome variables, we consider six indicators of patient health and utilization that could be related to an insufficient amount of PT: opioid prescriptions, pain management procedures, orthopedic surgeries, emergency department visits, inpatient stays, and skilled nursing stays. Given that the most common diagnosis for patients in our sample relates to muscle or joint pain, we might expect patients to substitute to opioids, pain management procedures, or orthopedic surgeries if PT did not successfully treat their pain. Patients also seek out PT to improve their strength and mobility—the top PT procedure code in our sample is for therapeutic exercises. Thus, an insufficient amount of PT could result in an injury or a fall that could result in an emergency department visit, an inpatient stay, or a skilled nursing stay.

Aside from opioid prescriptions, the outcomes are indicator variables that are measured within 12 months of a patient's last day of PT. We measure these outcomes starting after the last PT visit to capture patients seeking alternatives after PT has "failed," as opposed to contemporaneous utilization that could be endogenous to the PT's behavior—say, if a PT is co-located with a pain management practice. The opioid prescription outcome is an indicator variable that is measured between 12 and 24 months after the patient's last day of PT. We construct this measure using this time frame because the Part D prescription data for opioids is only available starting in 2006, so we cannot measure opioid prescriptions within 12 months for the 2005 sample. Appendix Section E provides further detail in how each outcome is constructed.

Before describing the results, it is important to consider what conclusions we can draw from our estimates, given our identification strategy and data. The therapy cap only restricts Medicare spending, so our estimates tell us the causal health effects of reducing *Medicare*- funded PT spending, rather than the effects of reducing *any* PT spending. If patients compensate for the reductions in Medicare PT spending by paying for the rest of their care out-of-pocket, then the therapy cap just represents a reduction in transfers, but not in actual PT utilization, for patients. In practice, patients likely do not fully compensate for these reductions, so the therapy cap should reduce overall PT spending. However, since our data is restricted to Medicare claims, we cannot measure how much a reduction in Medicare-funded PT spending passes through to overall PT spending. We can address this partially by repeating our analyses just on low-income patients, for whom we might expect more passthrough of the cap savings onto total spending. Furthermore, since we are using claims data, we cannot observe some dimensions of health that PT is arguably most relevant for, such as mobility and pain levels. Our analysis is not be able to speak to whether patient well-being worsens in these dimensions unless the effects are large enough to induce substitution to other forms of care or trigger events like ED visits.

Results Figure 4 shows the reduced form results from Equation 1. Panel (a) plots the relationship between first month of PT, interacted with an indicator for 2006, and PT spending in the subsequent 12 months. Consistent with the cap binding more for patients who start earlier in the year, there is a monotonic negative relationship between start month and 12-month PT spending—the reduction in PT spending due to the cap is much larger for patients who start earlier in the year relative to those who start later. Turning to our health outcomes, we would expect that if the cap-induced savings were poorly targeted, then there would similarly be a negative relationship between these outcomes and a patient's start month. Instead, this relationship is flat and all coefficients are statistically insignificant at the 5 percent level.



Figure 4: Reduced Form Spending and Health Outcomes

This figure plots the coefficient θ_f , which denotes the interaction between an indicator for 2006 and an indicator for month of first PT, from Equation 1. Panel (a) plots the coefficients on 12-month PT spending (\$). Panel (b) plots the coefficients on an indicator for pain management procedures, panel (c) plots the coefficients on an indicator for orthopedic surgery, panel (d) plots the coefficients on an indicator for emergency department visit, panel (e) plots the coefficients on an indicator for a hospital stay, panel (f) plots the coefficients on an indicator for a skilled nursing facility stay, and panel (g) plots the coefficients on an indicator for opioid prescriptions. All outcomes other than opioid prescriptions are measured within 12 months of the first PT session; opioid prescriptions are measured 12-24 months after the first PT session. Section E describes how the health outcome measures are defined. Sample is restricted to beneficiaries with more than \$200 in 12-month PT spending. Data: 20% Medicare Carrier, Outpatient, MEDPAR, and Part D claims.

In Figure 5, we plot the IV coefficients from estimating Equation 2. Given our relatively strong first stage,²¹ we are able to rule out out effect sizes larger 4 percentage points in both directions for all outcomes.



Figure 5: IV: Effect of PT Spending on Health Outcomes

This figure plots the coefficient β , which denotes the effects of an additional \$100 of PT on an indicator of each outcome, from Equation 2. All outcomes other than opioid prescriptions are measured within 12 months of the first PT session; opioid prescriptions are measured 12-24 months after the first PT session. Section E describes how the health outcome measures are constructed. Sample is restricted to beneficiaries with more than \$200 in 12-month PT spending. Data: 20% Medicare Carrier, Outpatient, MEDPAR, and Part D claims.

However, the lack of health impacts on the overall sample could be masking effects for particularly vulnerable patients. We assess this possibility in three ways. First, we focus our attention on two patient populations where the reduction in Medicare PT spending could be particularly impactful. We first consider low-income patients, who may be unwilling or unable to pay for additional PT care out-of-pocket in the face of a Medicare denial. If so, the reductions in Medicare PT spending from the cap should pass through into relatively larger reductions in *overall* spending, which is confirmed by the larger reduced form effect on 12-month PT spending in Figures B7 and B8. Figure B9 shows that for two types of low-income patients—dual-eligibles and patients receiving the Part D LIS—there is no evidence that the reduction in PT spending due to the cap worsened health outcomes.²²

²¹The Kleibergen-Paap Wald F-statistic for the first stage is 251.4.

²²While the reduced form results show a statistically significant and positive estimate on opioid prescrip-

Next, we stratify by a patient's initial health status, which is defined as their decile of predicted 2004-2005 PT spending constructed in Section 2.3. We estimate the IV in Equation 2 separately for each decile, where higher deciles indicate a patient is predicted to require more PT. We would expect to see negative effects to be concentrated mostly on relatively high-need patients. However, Figure B10 shows that even among the highest deciles of predicted PT spending, we find null effects on health outcomes.

Finally, we stratify by whether the patient recently received PT-related care prior to starting PT. As shown in Table A2, about five percent of patients had an orthopedic surgery and about a third of received a pain management procedure *prior* to starting PT. Thus, for these patients we may not expect that they would receive a second procedure or surgery after PT. Figure B11 stratifies the IV results by whether a patient reported having a pain procedure or orthopedic surgery within the last 6 months. Even among patients without a recent procedure or surgery, there is no evidence that reductions in PT spending led to increases in utilization of PT alternatives.

4 Characterizing Who is Screened Out by the Cap

Thus far, we have measured the extent to which the 2006 therapy cap affected PT spending and patient health overall. While the cap indeed reduced spending, it did not appear to worsen patient health, suggesting that the savings were well-targeted. In this section, we explore these targeting effects further by focusing on patients as they get close to cap and characterizing *who* the cap screens out and *why*. Given that the stated goal of the cap was to restrict medically unnecessary care, in the first set of analyses we compare patients of different levels of clinical need to determine if lower-need patients are more likely to be screened out than higher-need ones.

Our second set of analyses then examines *horizontal inequity*, conditional on need—that is, whether patients with the same level of PT need face different probabilities of being screened out by the cap in a systematic way. Because the cap introduced a documentation

tions for patients with the earliest start months, the IV estimate in Figure B9 for opioids is not statistically significant.

requirement, we pay particular attention to variation related to provider administrative capacity, which we proxy for using provider size. Furthermore, since provider size is also correlated with patient demographics in our context, we further explore whether differences across providers of different sizes translate into disparities across patient groups.

Figure 3 shows that the cap delivers savings through two distinct channels: by deterring attempts to bypass the cap and, for those who make an attempt, by denying them. The factors influencing whether a patient makes it past each stage differ—the deterrence decision reflects a joint decision to between the patient and their provider to continue care, while a denial depends on Medicare's evaluation of the documentation and enforcement of its medical necessity standard. We will thus consider the screening properties of the deterrence and the denial channels separately.

4.1 Methodology and Sample Construction

To characterize who is screened out by the cap, we focus on the weeks in which a patient approaches or makes an attempt to go over the cap and then examine the relationship between clinical need, whether they stop at the cap, and why they stop — deterrence or denial. We will isolate the targeting properties of the cap by looking at whether the slope of the relationship between need and deterrence or denial *changes* from 2005 to 2006. In particular, if the slope becomes steeper after the cap is implemented, this would indicate that the cap improved targeting on need through that particular channel.

Figure B12 panels (a) and (b) illustrate the intuition of this exercise. Each panel plots clinical need against the likelihood that a patient stops at the cap, separately for 2005 and 2006. For illustrative purposes, we do not differentiate between stopping due to deterrence or denial; when we actually implement this approach, we will evaluate each outcome separately. In 2005, the therapy cap is not in effect yet, so the "share stopping at cap" is simply the natural rate at which patients stop where the cap will eventually be located. In both panels, the downward-sloping relationship indicates that higher-need patients are at baseline less likely to stop at the cap. This relationship can be negative even before the cap is in place. This would be the case if, for example, even absent the cap, providers tend to give more care to sicker patients, making them more likely to stop at any arbitrary spending level for

relatively healthier patients.

Figure B12 panel (a) illustrates the case where the cap increases the share of patients who stop at the cap, but does *not* improve targeting on need. To see this, note that the likelihood of stopping at the cap increases uniformly at all levels of need such that the slope of this relationship remains unchanged. So while all patients were more likely to stop, high vs. low-need patients were not differentially affected. In contrast, in panel (b), while all patients are more likely to stop at the cap, the increase is *largest* for low-need patients. In other words, the cap tends to screen out lower-need patients, which implies an improvement in the targeting of spending on need. This intuition forms the basis for our empirical test of improved targeting: whether the slope of the relationship between stopping at the cap and patient need steepens between 2005 and 2006.

Figure B12 panel (c) then illustrates the intuition of our test for horizontal inequity on an indicator for dimension D (e.g., D could be an indicator for going to a large or small provider). Holding fixed patient need, patients with D = 1 in the pre-period, represented by the circles and dotted lines, are no more likely to stop at the cap than patients with D = 0. However, once the cap is in place, a gap emerges between patients with D = 1and D = 0—at the same level of need, patients with D = 1 are more likely to stop at the cap. We interpret this as indicating that the cap has introduced horizontal inequity along dimension D. Thus, our test our horizontal inequity is whether there is a level difference between patients with D = 0 and D = 1 that emerges in 2006 among patients with the same clinical need.

We first present our results as binscatters of the relationships between two separate outcomes at the cap—deterrence and denials—and our measure of patient need. When looking at deterrence, we focus on the sample of patients in the week they approach the cap—that is, the week before an actual attempt or, among those who do not attempt, the point where one more week is extrapolated to take them over the cap. For denials, we focus on patients in the week of their attempt to go over the cap—that is, the week where the cumulative paid amount before that week plus the billed amount for that week is projected to go over the cap. The relationships are residualized of attempt-week-times-year or approachweek-times-year fixed effects to account for secular time trends. We formally test for slope changes between 2005 and 2006 by estimating the following equation:

$$Y_i = \beta_1 X_i + \beta_2 X_i \times 1(Year_{y(i)} = 2006) + Week_{t(i)} + \varepsilon_i, \tag{3}$$

where Y_i is a dummy variable indicating either if the patient was deterred or denied, X_i is a measure of patient need, and $Week_{t(i)}$ is a fixed effect for the week-year that the patient approached or made an attempt.²³ β_1 is the slope of the relationship between X_i and the outcome Y_i in 2005, and β_2 captures how that slope changes in 2006. Thus, our test of whether targeting on need improves after the cap is introduced is whether β_2 is statistically significant and going in the same direction as β_1 , which would indicate a steepening of the slope.

Then to test for horizontal inequity conditional on need, we estimate the following equation:

$$Y_i = \beta_1 X_i + \beta_2 X_i \times 1(Year_{y(i)} = 2006) + \beta_3 D_i + \beta_4 D_i \times 1(Year_{y(i)} = 2006) + Week_{t(i)} + \nu_i, \quad (4)$$

where D_i is an indicator variable for whether a patient goes to a small provider or a patient demographic characteristic. β_3 captures whether patients with $D_i = 1$ are at baseline more or less likely to stop at the cap, and β_4 captures whether this changes once the cap is in place, after controlling for patient need. Thus, our test for whether the cap introduced horizontal inequity is whether β_4 is statistically significant.

We use the predicted 2004-2005 PT spending measure discussed in Section 2.3 as our main measure of patient need, X_i . Interpreting predicted pre-reform spending as a measure of underlying need requires assuming that a patient's relative position in the spending distribution is informative about their true latent need, and that this ranking is stable between 2005 and 2006. Thus even if the overall *level* of spending in 2004 and 2005 may have been

 $^{^{23}}$ As discussed in Section 2.3, a given patient can spend multiple weeks where they are labeled as approaching or attempting the cap. The former can happen because the weeks to the cap calculation is based on an extrapolation, so patient can spend multiple weeks looking like they are a week away from approaching the cap. The latter can happen if patients who are initially denied make multiple attempts across several weeks. We include observations for each approach and attempt week to properly account for week-year time trends, but the outcome variables are defined at the patient-level to capture if the patient *ever* attempts or *ever* makes it past the cap by the end of the year. To account for the possibility of multiple observations per patient, we cluster our standard errors at the provider-level.

considered too high by Medicare, our test is valid as long as *relatively* high-spend patients in 2004-2005 have greater true PT need than relatively low-spend ones.

One potential limitation of using a predicted outcome trained on healthcare claims as our measure of clinical need is that these prediction models are known to inherit biases and reproduce human judgment errors (Mullainathan and Obermeyer, 2017). For example, lowerincome populations may under-utilize care relative to their actual need because they are more sensitive to healthcare costs or due to provider bias. If prior utilization predicts greater future spending in the model, then it would underestimate true need for these populations. This is an issue inherent in any analysis using claims-based measures to proxy for patient need.

This potential for bias has several implications for the interpretation and validity of our results. For the targeting analyses that test for a change in slope, if predicted and true need were only weakly positively correlated or completely uncorrelated with each other, then this would bias our test *against* finding an improvement in targeting. Using a biased measure of need would undermine our analysis only if the bias is large enough such that the predicted and true need were sufficiently *negatively* correlated with each other, which is a fairly extreme case of model misspecification. In that case, what we would classify as an improvement in targeting (i.e., a steepening of the slope) would actually be indicative of worsened targeting. For the horizontal inequity analyses that test for a level change, this bias should generate a gap in outcomes between $D_i = 1$ patients and $D_i = 0$ patients in both the pre-reform and post-reform period. We need that the magnitude and distribution of this error with respect to X_i do not change over time in order for it to not affect our coefficient of interest β_4 , which is on the interaction between D_i and an indicator for $Year_{y(i)} = 2006$.

One way we can validate that our results are robust to these concerns is by repeating our analysis using patient age instead of predicted need in Figure B13 and Table A4. Age is plausibly correlated with true need and but not prone to claims-based measurement error. We find qualitatively similar results when using patient age as our measure of need.

4.2 Results: Targeting on Need

Figure 6 panel (a) considers the patients who approach the cap—patients for whom one more week of PT will take them over the cap—and plots the correlation between need and whether they eventually make an attempt. We interpret patients who approach but do not attempt as being "deterred."

Both before and after the cap, PTs are less likely to be deterred for patients with higher need, as indicated by the negative slope in 2005. Note that in the absence of the cap in 2005, "deterrence" is simply the natural rate at which patients stop care at the point where the cap will eventually be. Once the cap is in place in 2006, they are noticeably more likely to be deterred—on average, there is a 43 percent (7.9 percentage point) increase in the likelihood that a patient stops at the cap without an attempt from 2005 to 2006. However, the magnitude of this increase is similar across all levels of need, as reflected in the uniform upward shift in levels and the lack of a change in the slope. It does not appear that the cap leads to improved targeting on need in the deterrence stage, where patients and providers are deciding whether to make an attempt, conditional on approaching.

Figure 6 panel (b) then examines patients who make an attempt to go over the cap and considers how likely they are to be denied by Medicare. In contrast to the deterrence results in panel (a), here we see that Medicare's denials become more sensitive to patient need once the cap is in place. Prior to the cap, Medicare rarely denied PT claims, and the denial rate only has a slight negative correlation with need. Once the cap is introduced, Medicare became noticeably more likely to deny care, as shown in the shift up in denial rates at all levels of need. However, the magnitude of this increase differs by patient need and is substantially larger for low-need patients, which causes a steepening of the slope of this relationship.

We confirm these patterns formally in Table A3, which shows that the coefficient on "Predicted spending \times 2006" is insignificant for deterrence (column 1) but statistically significant and negative for denials conditional on attempt (column 3). Column 5 shows that the net effect on whether a patient eventually gets past the cap, conditional on *approach*, demonstrates an improvement in targeting, as the coefficient for denial conditional on approach on "Predicted spending \times 2006" is negative.



Figure 6: Correlation between Patient Need and Deterrence or Denial, 2005-2006

(a) Deterred vs. patient need,

(b) Denied vs. patient need, conditional

This figure plots the relationship between cap outcomes and patient need and provider size. "Deterred" is defined as the share of patients who approach the cap but do not attempt, and "denied" is defined as the share of patients who attempt but never make it past the cap. Panels (a) and (b) plot the relationship between log predicted 12-month PT spending and share deterred and denied in 2005 and 2006. Panels (c) and (d) plot the same relationship, split by provider size. Provider size is defined as the total number of Medicare beneficiaries who receive regular PT by that provider in 2006-2008, and a large firm is defined as being above-median. Section 2.3 and Appendix Section D discuss the construction of the predicted PT spending measure and Section 2 describes the sample definition in further detail. Data: 20% Medicare Carrier claims and Master Beneficiary Summary File.

These results show that while the overall savings from the cap are due to a mix of increased provider deterrence and Medicare denials, the improvements in *targeting* come only from Medicare's behavior. The extent to which providers pulled back care in response to the cap did not differ by a patient's *ex ante* need, indicating that the ordeal acts as a fairly blunt tool in terms of influencing care decisions. This blunt deterrence response is perhaps even more surprising since the ordeal itself— the cost of providing care as well as of producing documentation, minus the expected reimbursement—effectively *was* targeted. The likelihood of a cap denial is much higher for low-need patients, meaning the expected reimbursement (i.e., reimbursement net of denials) is lower for these patients. Instead, the uniform deterrence response across all levels of need suggests that the decision to continue past the cap was instead driven by factors uncorrelated with patient need. It could, for example, reflect the provider's opportunity cost of time spent on the visit, their willingness to bill patient in the event of a denial, or the patient's own ability to pay out of pocket if denied. However, it is important to note that providers' baseline care decisions were already targeted on need, as demonstrated by the negative slope between need and deterrence in 2005. The lack of further steepening in the slope in 2006 indicates that the cap itself did not induce any *additional* targeting in deterrence above and beyond that.

Given the potential concerns with using a claims-based prediction as an index of need, as additional robustness we repeat the analysis and estimate Equation 3 using patient age as our measure of need. The results are reported in Figure B13 panels (a) and (b), and Table A4 columns 1, 3, and 5. They show patterns consistent with the main results: the cap improves screening only in the denial stage.

4.3 **Results: Horizontal Inequity**

By Provider Size We next consider whether the cap "screens" on a particular characteristic which should be orthogonal to patient need: their provider's administrative capability. Given that the cap introduced a documentation requirement and was enforced through billing denials, we focus on differences by provider size since prior work has shown that larger providers tend to have an advantage in billing and documentation. Dunn et al. (2024) show that smaller providers incur higher billing costs, and League (2022) has shown that one way providers respond to increasing insurer denials is by consolidating into larger groups. Differences by provider size could arise in the deterrence stage if larger providers have lower documentation costs and are across the board more likely to make attempts, or they could manifest in the denial stage if larger providers have greater awareness of or compliance with billing requirements.

To test for disparities by provider size, Figure 6 panels (c) and (d) plot the same correlations as before between patient need and deterrence or denials, but split by whether a patient goes to a provider with above-median Medicare patient count—a "large provider"— or below-median—a "small provider."²⁴ Panel (c) shows that at a given level of need, there is no difference in deterrence across large or small providers in 2005 as well as 2006. In contrast, panel (d) shows that once the cap is introduced, patients who go to small providers are much more likely to be denied. The gap in denials between large and small providers exists even conditional on patient need, indicating that it cannot be attributed to differences in patient composition across different providers. It is instead likely to be driven by differences in *provider* behavior that influence the denial rate, independent of patient characteristics. The estimates from Equation 4 in Table A3 columns 2, 4, and 6 confirm these patterns formally.

By Patient Race and Income The difference in denials between large and small providers conditional on need is itself noteworthy because it suggests that the cap is screening on a characteristic unrelated to patient medical necessity, in contrast to the intention of the policy. However, these differences would have limited implications for horizontal inequity across *patients* if patient sorting across providers is largely random. But if different patient groups systematically see differently-sized providers, then screening on provider size is more troubling as it would generate horizontal inequity across these groups. Indeed, we find strong evidence of non-random patient sorting: patients of smaller providers, who are more likely to be screened out, tend to be lower-income and more likely to be minorities. Figure B14 plots the correlation of the relationship between patient characteristics and an indicator variable for going to a small provider, among patients who approach the cap in 2006. Patients who go to small providers tend to have lower ZIP income, are more likely to be dually-eligible for Medicaid or receive the Part D Low Income Subsidy (LIS), and are more likely to be non-white.

 $^{^{24}}$ The median patient count in our sample is 925 Medicare patients in total between 2006 and 2008.

These correlations between provider size and patient demographics thus motivate the analyses in Figure 7, which repeats the horizontal inequity analysis in Figure 6, splitting by patient demographic groups. Table A5 shows the estimates of the coefficient on predicted spending, predicted spending interacted with a 2006 year dummy, and a respective demographic indicator variable, as well as the demographic indicator interacted with a 2006 year dummy. Columns 1-3 show that among those who approach the cap in 2006, lower-income and minority patients are actually slightly less likely to be deterred by the cap, as indicated by the negative coefficient on the interaction between each demographic characteristic and the 2006 year indicator. In contrast, columns 4-6 show that among those who attempt, they are 19-38 percent more likely to be denied. On net, the effect in the denials channel outweighs the one on deterrence — among all those who approach the cap, non-white and lower-income patients are 8-11 percent more likely to be stopped by it, conditional on patient need (columns 7-9).
Figure 7: Correlation between Patient Need and Deterrence or Denial by Patient Demographics, 2005-2006



This figure plots the relationship between attempt outcomes and patient need and provider size. "Deterred" is defined as the share of patients who approach the cap but do not attempt, and "denied" is defined as the share of patients who attempt but never make it past the cap. Panels (a) and (b) plot the relationship between log predicted 12-month PT spending and share deterred and denied in 2005 and 2006, split by patient race. Panels (c) and (d) plot the same relationship, split by patient dual eligibility status for Medicaid in the respective year. Panels (e) and (f) plot the same relationship, split by whether the patient receives the Part D Low Income Subsidy in 2006. Section 2.3 and Appendix Section D discuss the construction of the predicted PT spending measure and Section 2 describes the sample definition in further detail. Data: 20% Medicare Carrier claims and Master Beneficiary Summary File.

The emergence of the gap by provider size in Figure 6 and the strong correlation between provider size and patient demographics in Figure B14 are consistent with patient sorting across large vs. small providers as being the key driver of the race and income gaps in Figure 7. However, an alternative explanation could be that these demographic gaps reflect across-group differences unrelated to the identity of a patient's provider. For example, the differences in deterrence could be due to low-income and minority patients being less healthy than their higher-income, white counterparts in a way that is observable to PTs but that is not captured in our predicted spending measures. The differences in denials could be driven by Medicare's denial algorithm inheriting some bias embedded in the data-generating process (Obermeyer et al., 2019). These should both generate demographic disparities both across- and within- providers. Thus, we test for this alternative explanation by estimating Equations 3 and 4 with the inclusion of provider-year fixed effects in Table A6. Once provider fixed effects are included, the disparities by income become statistically insignificant and the disparity by race shrinks substantially. This confirms that the disparities in Table 7 are largely driven by patient sorting across different providers — low-income and minority patients tend to go to providers who are less-adept at navigating the cap.²⁵

Table A7 calculates the differences in 2006 denial rates implied by the estimates in Tables A3 and A5 across provider size, patient race, and patient income, for the median patient who attempts or approaches the cap. Conditional on making an attempt, a patient with median predicted spending who goes to a small provider is 80 percent more likely to be denied than one who goes to a large provider. The minority-to-white gap in denials is 43 percent, the dual-to-non-dual-eligible gap is 23 percent, and the LIS-to-non-LIS gap is 21 percent.

Taken together with the targeting improvements from Section 4.2, these results demonstrate that screening with a soft spending limit introduces an equity-efficiency tradeoff. On the one hand, the cap improves targeting by medical necessity—patients with lower *ex ante* need see larger increases in denials than patients with higher need from the cap. On the other hand, the cap also introduces a substantial advantage for larger providers. Prior to the cap, two patients with similar medical need but who see differently-sized providers had similar

²⁵Furthermore, the model which includes provider fixed effects shows that the improvements in Medicare screening on denials are *not* the result of patient sorting. This can be seen by the negative and statistically significant estimates for the "Predicted staffing \times 2006" coefficients in Table A6, columns (5) and (9).

chances of continuing care; after the cap, the patient who sees the smaller provider has a much lower chance of continuing. As a result, patients who tend to see smaller providers, like non-white and lower-income populations, are more likely to be stopped at the cap. This "screening" on provider size appears to be driven by a factor unrelated to patient need. In the next section, we argue that it is driven by across-provider differences in administrative capacity, specifically in their knowledge of and compliance with the documentation requirement.

5 Drivers of Provider Size Advantage

The results in Section 4 demonstrate that there is substantial variation in denials across large and small providers, which in turn translates into differences across patient race and income groups that is orthogonal to patient clinical need. In this section, we show that these across-provider differences are the result of variation in their compliance with the cap's documentation requirement. Furthermore, the reason large providers tend to do better than smaller ones is because of they have more opportunities for learning-by-doing in compliance with the cap requirements.

5.1 The Role of Documentation

While documentation review was not conducted for every approved exception to the cap, CMS emphasized the importance of always having documentation *available* for review in its communications about the therapy cap (CMS, 2006b). Providers were instructed to add a modifier code (the "KX modifier") to their claims to attest that documentation was available. We cannot observe the documentation itself or verify that all providers who claimed to have it actually did, but we interpret the use of the modifier code as an indicator of baseline provider awareness about the requirement to have documentation on hand to support their request.

Documentation use is strongly associated with denials both in the cross-section and in over time. Table 1 first explores the cross-sectional variation, and reports the estimates of a regression between an indicator for using documentation on an attempt in 2006 and the likelihood that the attempt is approved. There is a strong positive correlation between the two, and the inclusion of documentation approximately doubles the likelihood that an attempt is approved. The magnitude of the coefficient is stable even with the inclusion of patient demographics, patient predicted health, and provider size. Including the indicator for having documentation into the regression model increases the R^2 by 2-4-fold.

We next examine the variation in documentation use and denials over time. Figure 8 panel (a) plots the share of attempts with documentation and the share of attempts approved between 2004 and 2009. Prior to the introduction of the cap, documentation use is mechanically zero, since the modifier code was not introduced yet. Strikingly, once the cap is in place, the claim approval rate moves in lock-step with the use of documentation. This co-movement strongly suggests that changes in documentation use are driving the increase in approvals over time.²⁶

Documentation use (and therefore the approval rate on attempts) is also highly correlated with provider size: large providers are much more likely to use documentation on attempts than small providers. The binscatters in Figure 8 panels (b) and (c) show that provider size is strongly positively associated with both documentation use and approvals. Providers in the top quartile of size, who see about 1500 Medicare patients from 2006-2008, are 22 percent (11.8 percentage points) more likely to use documentation on their attempts in 2006 than those in the bottom quartile, who see 575 patients. This translates into a sizeable differences in approvals by size—attempts by providers in the top quartile of size are 48 percent (16.1 percentage points) more likely to be approved than those in the bottom quartile.

 $^{^{26}}$ In Appendix Section D, we use a machine-learning model to investigate whether the change in approval rates over time could reflect changes in claim-level strategic *billing* behavior that is separate from documentation, such as upcoding.

	(1)	(2)	(3)	(4)
		Appi	roved	
Has documentation	22.41***	21.89***	22.49***	20.85***
	(0.887)	(0.905)	(0.880)	(0.908)
N	51049	43620	54455	43620
\mathbb{R}^2 , no modifier code	.018	.019	.024	.032
R^2 , with modifier code	.068	.067	.074	.075
Outcome mean	46.88	47.02	46.25	47.02
Patient demographics	Х	Х		Х
Patient health		Х		Х
Provider size			Х	Х

Table 1: Regression of Documentation on Approvals on Cap Attempts, 2006

This table reports the relationship between approvals and documentation on attempts in 2006. The outcome variable is the share of claims where the patient faces no denials in the week that they make an attempt, where the attempt week is defined as in Section 2.3. "Has documentation" is an indicator variable for whether any claim submitted that week has a KX modifier code. The table reports two R^2 values: one is for the regression which includes the "Has documentation" variable in addition to all the controls listed below, and the other is for a regression which includes only the control variables. The patient demographic control variables are age, race, and sex. Patient health controls are predicted 12-month PT spending, in-office spending in 6 months prior to first PT, total inpatient spending last year, SNF spending last year, total Part B spending last year, and total imaging spending last year. Provider size is measured as the total number of attempts the patient's provider made that year. All specifications include week-year fixed effects for the calendar week(s) that a beneficiary attempts to go over the cap. Standard errors (in parentheses) are clustered at the provider level. Data: 20% Medicare Carrier claims and Master Beneficiary Summary File.



Figure 8: Documentation, Approvals, and Provider Size

(a) Time Series, 2006-2008

(b) Documentation vs. Provider Size



(c) Approvals vs. Provider Size

Panel (a) plots the share of attempts with at least one claim with documentation ("Share with documentation") and without any denials ("Share approved") in 2006-2008. "Share with documentation" is the share of attempts with at least one claim using the KX modifier code. Panels (b) and (c) plot the relationship between provider size and approvals and documentation use (respectively) on attempts in 2006. Attempts are defined as in Section 2.3: weeks in which a patient's cumulative paid amount is below the cap at the beginning of the week, and their cumulative paid or billed amount is above the cap at the end of the week. Provider size is measured as a TIN-state's 2006-2008 Medicare patient count. Data: 20% Medicare Carrier claims.

5.2 Link between Provider Size and Documentation: Learning

Decomposition of Provider Size Advantage We next explore *why* larger providers have better compliance with the documentation requirement and thus fare better in the face of the cap. Documentation practices could vary across providers of different sizes for a variety of reasons. On the one hand, large and small providers may just vary in their behavior at baseline: they could have made differing initial investments in technology like electronic health records or have more knowledgeable billing staff. The correlation could also reflect an omitted factor that is correlated with both provider size and documentation use, like ownership status or geographic location.

On the other hand, the increase in documentation and approvals throughout the first year in Figure 8 suggests an alternative mechanism for the size advantage: learning. If providers learn through their experience with the cap — "learning-by-doing"—then large providers naturally derive an advantage from the ability to move up the learning curve more quickly. Figure 9 panel (a) shows that small providers—those below the median accumulate experience with the cap much more slowly than large providers. By the end of the first year, approximately 60 percent of small providers still have less than 10 patients worth of experience with the cap, whereas the majority of large providers have more than 20 patients worth of experience.

We next decompose the variation in documentation use into learning-by-doing and persistent size-based advantages. We estimate a linear model of an indicator for documentation use in the week of an attempt, Y_i , for patient *i* receiving care from provider *j* that is attempting to bill above the cap in week *t*:

$$Y_{i} = \sum_{\substack{e(j(i),t(i))=1\\\text{Learning-By-Doing}}}^{50} \kappa^{e} + \underbrace{Week_{t(i)}}_{\text{Industry-Wide Trends}} + \underbrace{\alpha_{j(i)}}_{\text{Provider FE}} + \underbrace{\beta X_{i}}_{\text{Patient Controls}} + \varepsilon_{i}$$
(5)

The regression decomposes the variation into three key components of interest. The first is κ^e , which captures the average amount of learning-by-doing by a provider who has experienced e(j(i), t(i)) previous patients approaching the cap. The second is calendar week fixed effect $Week_{t(i)}$, which captures any industry-wide trends affecting all providers. This could encompass, for example, widespread dissemination of information about documentation requirements across all providers by CMS or industry groups. The last is $\alpha_{j(i)}$, a provider fixed effect that captures any persistent provider-level difference in documentation use. This captures baseline variation in provider aptitude for billing that may result from investments or other provider characteristics. We also control for X_i , the predicted need of the patient associated with the attempt.

Figure 9 panels (b)-(d) present the estimates from Equation 5. Panel (b) shows the estimated economies-of-learning curve and plots κ^e for $e \in [1, 50]$; the estimates suggest a significant amount of learning-by-doing, especially earlier on. The implied advantage from being farther down the learning curve can be quite large, with a veteran provider with more than 50 patients worth of experience having a 20 percentage point advantage over a complete novice. Panel (c) shows our estimates of industry-wide trends $Week_t$ and shows significant industry-wide increases in approvals that occurred throughout the first year. Finally, panel (d) plots a binscatter of provider fixed effects α_j against provider size. The lack of a relationship between the fixed effects and provider size indicates that, after accounting for differences in experience, larger providers do not have a baseline advantage over smaller ones in compliance with the documentation requirement. Taken together, the results indicate that the provider size advantage in documentation stems primarily from learning-by-doing, which mechanically accumulates faster for larger providers.²⁷

²⁷We repeat this analysis in Figure B15 where the outcome is approvals of attempts instead of documentation use. Here, we also find a learning curve consistent with learning-by-doing, but the provider fixed effects indicate that even after accounting for learning-by-doing about documentation use, large providers have some fixed advantage on approvals. This suggests that larger providers may have some baseline advantage with approvals that extends beyond basic documentation use. For example, they could have access to technology to produce richer documentation of medical necessity or more knowledgeable billing staff.



Figure 9: Decomposition of Size Advantage on Documentation Use

This figure characterizes differences in documentation use on cap attempts by provider size and experience in 2006-2008 using data and regression estimates from Equation (5). Panel (a) plots the share of large and small providers who achieve a cumulative number of cap attempts over time, where large providers are defined as having abovemedian patient count in 2006-2008. Panel (b) plots the coefficient on the provider's cumulative number of prior attempts on approvals. Panel (c) provides estimates of weekly industry-wide trends. Panel (d) plots the provider fixed effects against provider patient count. Section 2.3 describes the definition of cap attempts in further detail. Data: 20% Medicare Carrier claims.

Sharp Evidence of Learning-By-Doing As a final piece of evidence on provider learningby-doing about documentation, we look at how provider behavior evolves around events which should be associated with learning. For each provider, we identify the point at which they seem to "learn" how to avoid a denial by looking for the first time that the provider successfully reverses a denial on a previous attempt with a patient. In particular, we look for the weeks in which a provider makes an approved attempt *after* receiving a denial in an attempt with the same patient in a prior week.

We use a stacked event study method, following Cengiz et al. (2019). Each stack is centered around a given provider's learning event with a focal patient — the "treatment" — and we look at that provider's five attempt-weeks before and after their event, with patients *other* than the focal one. We use attempt-weeks instead of calendar weeks so as to only capture weeks in which the provider makes at least one attempt. These observations form the treated group within each stack. The control group within each stack are "clean controls" comprised of *other* providers who make attempts in the same calendar week as the learning event, but who do not have a change in treatment status in the 5 attempt-weeks before or after the focal week. In other words, they are either not treated in the entire period (because they are not-yet-treated or they are never-treated) or treated in the entire period. A given attempt can appear in the control group for multiple stacks, but can only be the learning event that defines treatment in one stack.

The regression specification for patient i's attempt within stacked group g is:

$$Y_{igt} = \sum_{a(i,g)=-5}^{5} 1(RelWeek_{a(i,g)} = a) \times Treat_{j(i),g} + \underbrace{\alpha_{j(i),g}}_{\text{Provider-Group FE}} + \underbrace{RelWeek_{a(i,g),g}}_{\text{Week Relative to}} + \underbrace{Week_{t,g}}_{\text{Calendar Week-}} + \underbrace{\beta X_i}_{\text{Patient Controls}} + \varepsilon_{igt}.$$
(6)

The outcome of interest, Y_{igt} is either documentation use for attempts or approvals on attempts in calendar week t. $Treat_{j,g}$ is an indicator variable for whether the attempt with patient i is associated with provider j's "learning event" and $\alpha_{j,g}$ is a provider-group fixed effect. $RelWeek_{a,g}$ is the attempt-week (between -5 and 5) relative to the learning event that defines group g, and $Week_{t,g}$ is the calendar week associated with patient i's attempt (interacted with group g). These can be separately identified because not all providers make an attempt in every week. We also control for X_i , the predicted need associated with the patient. Results are clustered at the provider-level to account for repeated observations across stacks.

Figure 10 plots the coefficients from Equation 6: panel (a) shows the results for documentation use during the attempt-week, and panel (b) shows the results for whether the attempt was approved (i.e., not denied). Both sets of results show that prior to the learning event, treated providers were not on a differential trend relative to control providers. However, after the provider corrects their first denial, they are consistently more likely to include documentation on future attempts and more likely to be approved. These results suggest sharp learning-by-doing within a provider—once they successfully reverse their first denial, they change their billing behavior to ensure future attempts are approved.





This figure plots the coefficients from the stacked event study specification around a provider's "learning event" in Equation 6. The outcome variables are (a) the share of attempts with documentation and (b) the share of attempts approved (i.e., not denied). The learning event is defined as a week in which a provider makes an approved attempt with a patient after receiving a denied attempt in a previous week with the same patient in a prior week. The sample is of attempts with patients *other* than the one associated with the learning event. An attempt-week is a week in which the provider made at least one attempt, as defined in Section 2.3. Each learning event is grouped and compared to other providers who also made attempts in the same calendar week, and the groups are stacked together (Cengiz et al., 2019). The specification includes controls for patient predicted PT spending, provider-group, and week-group fixed effects, and is clustered at the provider and group levels. Section 2.3 describes the definition of cap attempts in further detail. Data: 20% Medicare Carrier claims.

Discussion We have shown that patients of large and small providers face substantially different approval rates in their attempts to go over the cap in the first year. Furthermore, our decomposition and event study results demonstrate that this size advantage is the result of larger providers learning faster about correct documentation use. This implies that in the long-run, small providers' approval rates should approach that of large providers. Thus, while this short-run horizontal inequity across providers is likely undesirable, the longer-run implications of the size advantage depend on whether the policymaker's desired approval rate is closer to that of small or large providers—that is, whether their eventual goal is to for exceptions to the cap to be relatively rare or frequent.²⁸

If the goal is to maximize savings and for cap approvals to be relatively rare—as they are for small providers—then the fact that there is rapid provider learning means that the cap's savings and efficacy will diminish quickly over time. Consistent with this, Figure B17 shows that the excess and missing mass around the cap decrease substantially from 2006 to 2008, meaning the overall cap savings are diminishing over time. One way to mitigate this is to make Medicare's approval process less deterministic, and thus more difficult to "learn" about. Instead of having approvals depend mostly on whether the provider simply indicates that they have documentation, approvals could be tied more directly to the content of the documentation. The tradeoff here would be that in-depth documentation reviews are costly for both Medicare and providers.²⁹

However, if the goal is for the cap to be fairly "soft", meaning approvals are relatively frequent—as they are for large providers—then in order to get ahead of horizontal inequity by size, Medicare will interventions that level the playing field up front for smaller providers. Rather than letting providers learn through cumulative experience, Medicare could instead provide targeted provider education about billing rules or subsidize related technology—both of which are tactics it has used before (CMS, 2014, 2024b).

 $^{^{28}}$ In Figure B16, we re-run the health analysis from Section 3.2, split by patients of above- or belowmedian-sized providers. We find that there is no detectable health effect from cap-induced savings in either type. This suggests that making cap approvals relatively rare, as is the case for small providers, would not cut medically necessary spending.

²⁹For example, Medicare's Recovery Audit Contractor program contracted with clinicians to conduct indepth documentation reviews of medical necessity (Shi, 2024).

6 Conclusion

This paper studies how soft spending limits screen by examining a spending limit imposed on healthcare providers by Medicare for physical therapy. We find that the soft "therapy cap" reduced overall spending, though much less than when compared to a hard cap that limits all spending above the cap. We look for direct effects of these savings on patient health using an identification strategy that leverages the annual nature of the cap, which makes the cap more binding for patients who start PT earlier in the year. We find precise null effects of cap-induced spending reductions on PT substitutes like opioid use, pain procedures, and orthopedic surgeries, as well as on ED visits and inpatient or skilled nursing stays.

Using a novel feature of healthcare claims data which allows us to observe both successful and unsuccessful attempts to bypass the cap, we then characterize which patients are stopped at the cap and why. This allows us to differentiate between patients stopping due to the deterrence channel—those who stop just short of the limit without making an attempt—and the denial channel—those who make an attempt but are not approved by Medicare. We find that both channels contribute to the overall savings, with denials contributing slightly more half. We then assess the targeting effects of each channel by comparing the correlation between receiving care past the cap and predicted patient need. While spending decreases for patients across the board, the reductions are largest for lower-need patients. Thus, the cap improves the targeting of spending. Once we differentiate between the deterrence and denials channels, we find that the targeting improvements stem entirely from Medicare's actions in through the denials channel.

Furthermore, we find that the cap introduces inequities which did not exist before: conditional on need, patients who go to small providers were much more likely to be stopped by the cap, as their providers face a much higher Medicare denial rate. As lower-income and minority patients tend to see smaller providers, these differences across providers translate into large gaps in denials by race and income among patients with similar levels of observable need.

We find evidence that these differences are driven by heterogeneous administrative capability across providers. All else equal, attempts to go over the cap are much more likely if providers indicate they have documentation of need, and larger providers are much more likely to do so. We decompose the large provider advantage into the components driven by learning-by-doing with experience, secular changes over time, and fixed provider-specific effects. We find strong evidence of learning-by-doing: providers become more successful at navigating the cap with experience. This lends larger providers, who gain experience faster, a mechanical advantage. We also provide sharp evidence in an event study design showing that providers learn from billing mistakes—after the first time they reverse a cap denial, providers are much more likely to use documentation and receive an approval on subsequent attempts.

Taken together, our findings suggest that soft spending limits work, but with important tradeoffs. On the one hand, imposing such a limit does indeed improve the efficiency of spending—spending limits allow principals to strike a balance between reducing overall spending while still allowing high-value spending to occur. In our setting, these improvements in targeting stem entirely from additional regulator scrutiny rather than improved "self-screening" by providers. On the other hand, they also introduce perhaps undesirable horizontal inequity due to the associated paperwork requirements. In the case of the therapy cap, we find that employing a documentation requirement gives larger providers a substantial advantage simply because they have more opportunities to learn how to comply with the policy. The result is an efficiency-equity tradeoff: while these policies improve the targeting of spending, they also effectively screen on agents' administrative ability, leading to horizontal inequity.

References

- Alatas, Vivi, Abhijit Banerjee, Rema Hanna, Benjamin A. Olken, Ririn Purnamasari, and Matthew Wai-Poi, "Self-Targeting: Evidence from a Field Experiment in Indonesia," *Journal of Political Economy*, April 2016, 124 (2), 371–427. Publisher: The University of Chicago Press.
- Amico, Peter, Gregory C. Pope, Poonam Pardasaney, Ben Silver, Jill A. Dever, Ann Meadow, and Pamela West, "Refinements of the Medicare Outpatient Therapy Annual Expenditure Limit Policy," *Physical Therapy*, December 2015, 95 (12), 1638–1649.
- **APTA**, "MedPAC Releases Revised Medicare 'Payment Basics'," Technical Report October 2020.
- Brot-Goldberg, Zarek C., Samantha Burn, Timothy Layton, and Boris Vabson, "Rationing Medicine Through Bureaucracy: Authorization Restrictions in Medicare," January 2023.
- Cengiz, Doruk, Arindrajit Dube, Attila Lindner, and Ben Zipperer, "The Effect of Minimum Wages on Low-Wage Jobs^{*}," *The Quarterly Journal of Economics*, August 2019, 134 (3), 1405–1454.
- CFR, "Initial Determinations," 2009.
- _, "Approval of the justification," 2020.
- **CMS**, "Outpatient Therapy Caps: Exceptions Process Required by the DRA," February 2006.
- _ , "Use of the KX Modifier on Claims Submitted to the Fiscal Intermediary When Some Services Exceed the Therapy Caps | Guidance Portal," HHS-0938-2006-F-8208, Centers for Medicare and Medicaid June 2006.
- _ , "Therapy Caps and Advance Beneficiary Notice of Noncoverage (ABN), Form CMS-R-131, FAQs April 2013," Technical Report, Centers for Medicare and Medicaid April 2013.
- _ , "An Introduction to: Medicare EHR Incentive Program for Eligible Professionals," Technical Report April 2014.
- _, "Outpatient Therapy Services and Advance Beneficiary Notice of Noncoverage (ABN), Form CMS-R-131, August 2018," August 2018.
- _, "Outpatient Physical and Occupational Therapy Services (L33631)," Technical Report, Centers for Medicare and Medicaid 2020.
- _ , "Medicare Benefit Policy Manual. Chapter 15 Covered Medical and Other Health Services," Technical Report, Centers for Medicare and Medicaid 2024.

_ , "Targeted Probe and Educate | CMS," Technical Report, Centers for Medicare and Medicaid September 2024.

- Currie, Janet, "The take-up of social benefits," in "Public Policy and the Distribution of Income," Russell Sage Foundation, 2006, pp. 80–148.
- **Deshpande, Manasi and Yue Li**, "Who Is Screened Out? Application Costs and the Targeting of Disability Programs," *American Economic Journal: Economic Policy*, November 2019, *11* (4), 213–248.
- Diamond, Rebecca and Petra Persson, "The Long-term Consequences of Teacher Discretion in Grading of High-stakes Tests," April 2016.
- Dillender, Marcus, "What happens when the insurer can say no? Assessing prior authorization as a tool to prevent high-risk prescriptions and to lower costs," *Journal of Public Economics*, September 2018, 165, 170–200.
- **DOL**, "How Do I File for Unemployment Insurance?," Technical Report, U.S. Department of Labor 2025.
- Dunn, Abe, Joshua D Gottlieb, Adam Hale Shapiro, Daniel J Sonnenstuhl, and Pietro Tebaldi, "A Denial a Day Keeps the Doctor Away*," The Quarterly Journal of Economics, February 2024, 139 (1), 187–233.
- **DynCorp**, "Study and Report on Outpatient Therapy Utilization: Physical Therapy, Occupational Therapy, and Speech-Language Pathology Services Billed to Medicare Part B in All Settings in 1998, 1999, and 2000," September 2002.
- Eliason, Paul, Riley League, Jetson Leder-Luis, Ryan McDevitt, and James Roberts, "Ambulance Taxis: The Impact of Regulation and Litigation on Health Care Fraud," *Journal of Political Economy*, November 2024. Publisher: The University of Chicago Press.
- Evans, William N., Shawna Kolka, James X. Sullivan, and Patrick S. Turner, "Fighting Poverty One Family at a Time: Experimental Evidence from an Intervention with Holistic, Individualized, Wrap-Around Services," *American Economic Journal: Economic Policy*, 2024.
- **FEMA**, "Insurance Documentation Required During FEMA Application Process | FEMA.gov," Technical Report, Federal Emergency Management Agency March 2021.
- Finkelstein, Amy and Matthew J Notowidigdo, "Take-Up and Targeting: Experimental Evidence from SNAP*," *The Quarterly Journal of Economics*, August 2019, 134 (3), 1505–1556.
- Goldstein, Amy, "Medicare Cutbacks Prove Painful to Some," Washington Post, May 1999.

- Howard, David H. and Ian McCarthy, "Deterrence effects of antifraud and abuse enforcement in health care," *Journal of Health Economics*, January 2021, 75, 102405.
- Hoynes, Hilary W., Nicole Maestas, and Alexander Strand, "Legal Representation in Disability Claims," March 2022.
- Ida, Takanori, Takunori Ishihara, Koichiro Ito, Daido Kido, Toru Kitagawa, Shosei Sakaguchi, and Shusaku Sasaki, "Choosing Who Chooses: Selection-Driven Targeting in Energy Rebate Programs," September 2022.
- **IRS**, "SOI Tax Stats Individual income tax statistics ZIP Code data (SOI) | Internal Revenue Service," Technical Report, Internal Revenue Service 2025.
- Kleven, Henrik Jacobsen, "Bunching," Annual Review of Economics, 2016, 8 (1), 435–464. _eprint: https://doi.org/10.1146/annurev-economics-080315-015234.
- and Wojciech Kopczuk, "Transfer Program Complexity and the Take-Up of Social Benefits," American Economic Journal: Economic Policy, February 2011, 3 (1), 54–90.
- Kopczuk, Wojciech and Cristian Pop-Eleches, "Electronic filing, tax preparers and participation in the Earned Income Tax Credit," *Journal of Public Economics*, August 2007, *91* (7), 1351–1367.
- League, Riley, "Administrative Burden and Consolidation in Health Care: Evidence from Medicare Contractor Transitions," 2022.
- Leder-Luis, Jetson, "Can Whistleblowers Root Out Public Expenditure Fraud? Evidence from Medicare," *Review of Economics and Statistics*, 2023, *Forthcoming*, 60.
- Lieber, Ethan M. J. and Lee M. Lockwood, "Targeting with In-Kind Transfers: Evidence from Medicaid Home Care," *American Economic Review*, April 2019, 109 (4), 1461–1485.
- Luthra, Shefali, "Tucked into the budget deal, long-awaited gifts to some health-care providers," *Washington Post*, March 2018.
- Macambira, Danil, Michael Geruso, Anthony Lollo, Chima D. Ndumele, and Jacob Wallace, "The Private Provision of Public Services: Evidence from Random Assignment in Medicaid," August 2022.
- Mullainathan, Sendhil and Ziad Obermeyer, "Does Machine Learning Automate Moral Hazard and Error?," American Economic Review, May 2017, 107 (5), 476–480.
- Nichols, Albert L. and Richard J. Zeckhauser, "Targeting Transfers through Restrictions on Recipients," *The American Economic Review*, 1982, 72 (2), 372–377. Publisher: American Economic Association.

- **Obermeyer, Ziad, Brian Powers, Christine Vogeli, and Sendhil Mullainathan**, "Dissecting racial bias in an algorithm used to manage the health of populations," *Science*, October 2019, *366* (6464), 447–453. Publisher: American Association for the Advancement of Science.
- **OIG**, "A South Texas Physical Therapist Claimed Unallowable Medicare Part B Reimbursement for Outpatient Physical Therapy Services," Technical Report A-06-14-00064, Office of the Inspector General June 2016.
- ____, "Fox Rehabilitation Claimed Unallowable Medicare Reimbursement for Outpatient Therapy Services," Technical Report A-02-16-01004, Office of the Inspector General August 2017.
- ____, "Many Medicare Claims for Outpatient Physical Therapy Services Did Not Comply With Medicare Requirements," Technical Report A-05-14-00041, Office of the Inspector General March 2018.
- O'Malley, A. James, Thomas A. Bubolz, and Jonathan S. Skinner, "The diffusion of health care fraud: A bipartite network analysis," *Social Science & Medicine (1982)*, June 2023, 327, 115927.
- Shepard, Mark and Myles Wagner, "Do Ordeals Work for Selection Markets? Evidence from Health Insurance Auto-Enrollment," August 2024.
- Shi, Maggie, "Monitoring for Waste: Evidence from Medicare Audits*," The Quarterly Journal of Economics, May 2024, 139 (2), 993–1049.
- **SSA**, "Disability Evaluation Under Social Security," Technical Report 64-039, Social Security Administration 2008.
- US GAO, "Early Resolution of Overcharges for Therapy in Nursing Homes Is Unlikely," Technical Report GAO/HEHS-96-145 August 1996.
- WebPT, "CMS Final Rule: The Countdown to 2024 Begins," 2024.
- Woodward, Susan E. and Robert E. Hall, "Diagnosing Consumer Confusion and Suboptimal Shopping Effort: Theory and Mortgage-Market Evidence," *American Economic Review*, December 2012, 102 (7), 3249–3276.
- **Zeckhauser, Richard**, "Strategic sorting: the role of ordeals in health care," *Economics* & *Philosophy*, March 2021, 37 (1), 64–81.
- Zwick, Eric, "The Costs of Corporate Tax Complexity," American Economic Journal: Economic Policy, May 2021, 13 (2), 467–500.

A Appendix Tables

HCPCS	Description	N. Lines (1000s)	Share of Beneficiaries $(\%)$
97110	Therapeutic exercises	3065.8	20.3
97140	Manual therapy	1657.7	12.8
G0283	Elec stimulation other than wound	875.5	7.4
97035	Ultrasound therapy	843.5	8.3
97112	Neuromuscular reeducation	651.9	5.4
97530	Therapeutic activities	640.7	5.5
97032	Electrical stimulation	415.5	3.5
97001	PT evaluation	279.8	18.6
97124	Massage therapy	226.8	2.0
97116	Gait training therapy	184.2	1.8
97010	Hot or cold packs the rapy	130.9	1.3
97012	Mechanical traction therapy	102.1	1.1
97150	Group therapeutic procedures	85.2	1.0
97113	Aquatic therapy/exercises	84.8	0.7
97535	Self care management training	76.2	1.5

Table A1: Highest-Volume In-Office Physical Therapy Procedures

This table reports the 15 highest volume PT procedures in the 2006 20% Carrier sample. "N. Lines" is the number of lines billed for each procedure code (HCPCS). "Share of beneficiaries" is the share of beneficiaries in the bunching sample that ever receive that procedure in that year. Data: 20% Medicare Carrier claims.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Bunchin	g Sample	Health Sample		Approa	ch Sample	Attempt Sample	
	2005	2006	2005	2006	2005	2006	2005	2006
Demographics								
Age	73.0	73.1	72.3	72.4	73.3	73.5	73.2	73.4
White	0.88	0.88	0.90	0.90	0.87	0.86	0.87	0.86
Female	0.65	0.65	0.66	0.66	0.65	0.65	0.65	0.65
Urban	0.85	0.85	0.84	0.84	0.86	0.87	0.87	0.88
ZIP code income $(\$)$	67,699	68,735	67,166	68,131	70,668	71,508	71,216	71,910
Part D Low Income Subsidy	0.15	0.15	0.15	0.14	0.17	0.18	0.17	0.19
Dual eligible	0.14	0.15	0.14	0.22	0.16	0.17	0.16	0.18
Prior utilization and spending								
Any hospital stay last 6m	0.22	0.23	0.24	0.22	0.26	0.27	0.27	0.28
PT-related surgery last 6m	0.04	0.05	0.12	0.11	0.05	0.05	0.05	0.05
Pain procedure last 6m	0.31	0.33	0.35	0.36	0.32	0.34	0.32	0.34
In-office spending last 6m (\$)	2088	2249	2313	2320	2538	2771	2620	2912
Predicted 12-month PT spending (\$)			1450	1473	1618	1703	1649	1760
PT utilization/spending								
Number of visits	12.65	12.73	11.73	10.74	23.03	20.63	25.17	23.66
Number of weeks of PT	7.68	7.85	7.08	6.60	13.42	12.36	14.62	14.17
Total PT spending (\$)	1161	1114	1035	890	2426	2030	2688	2345
Average spending per visit (\$)	105.17	102.73	99.33	95.74	127.17	119.54	128.69	118.89
Week start PT	22	21	26	26	14	14	14	14
Observations								
Number of beneficiaries	125142	128841	62063	71019	38209	35923	31082	25126
Number of providers	13222	13084	9242	9520	8928	8784	8164	7354

Table A2: Summary Statistics

This table reports summary statistics for the analysis samples in 2005 and 2006. Columns 1 and 2 show the summary statistics for the bunching analysis sample in Section 3.1, columns 3 and 4 report statistics for the health analysis sample in Section 3.2, and columns 5-8 report statistics for the screening samples ("Approach" and "Attempt") in Section 4. The predicted 12-month PT spending measure is calculated only for the health analysis and screening samples. Spending measures are in terms of Medicare spending and do not include patient coinsurance. Receipt of the Part D Low Income Subsidy is measured in 2006 as the program was not available in 2005. Data: 20% Medicare Carrier claims, Master Beneficiary files, MEDPAR, and Outpatient files.

	(1)	(2)	(3)	(4)	(5)	(6)	
	Outcome: Deterred			Outcome	e: Denied		
	Cond. on	Approach	Cond. on	Attempt	Cond. on	Approach	
Predicted Spending	-16.9***	-16.9***	-1.54***	-1.54***	-20.6***	-20.6***	
	(.702)	(.703)	(.176)	(.177)	(.797)	(.798)	
Predicted Spending \times 2006	533	453	-6.22***	-6.86***	-3.29**	-3.72**	
	(1.03)	(1.03)	(1.34)	(1.32)	(1.49)	(1.48)	
Small provider		741		.154		371	
		(.513)		(.145)		(.612)	
Small provider \times 2006		-1.28		12.7***		9.72***	
		(.815)		(1.06)		(1.18)	
Outcome mean, 2005	18.4	18.4	1.0	1.0	21.1	21.1	
Outcome mean, 2006	26.3	26.3	23.1	23.1	52.2	52.2	
Week-year FE	Х	Х	Х	Х	Х	Х	
Cluster	Provider	Provider	Provider	Provider	Provider	Provider	
N. Providers	11911	11911	11064	11064	11911	11911	
N. Patients	70518	70518	53560	53560	70518	70518	
N. Observations	80532	80532	116360	116360	80532	80532	

Table A3: Regression Results on Screening and Differences by Provider Size using Predicted Spending, 2005-2006

* p < .10, ** p < .05, *** p < .01. This table presents the coefficients from estimating Equation 3 with log predicted PT spending as X_i (columns 1, 3, and 5) and Equation 4 with log predicted PT spending as X_i and an indicator for whether a patient goes to an above-median-size provider as D_i (columns 2, 4, and 6). The coefficient β_1 is the estimate for "Predicted spending," the coefficient for β_2 is the estimate for "Predicted spending $\times 2006$ ", the coefficient for β_3 is the estimate for "Small provider" and β_4 is the estimate for "Small provider $\times 2006$." All specifications include week-year fixed effects for the calendar week(s) that a beneficiary attempts or approaches the cap. The regression is clustered at the provider (TIN-state) level. Provider size is measured as a TIN-state's 2006-2008 Medicare patient count and "small provider" denotes a provider with below-median patient count. Columns 1-2 and 5-6 restrict to the sample of 2005-2006 patients who ever approach the cap and columns 3-4 restrict to the sample of patients who ever make an attempt to go past the cap, as defined in Section 2.3. Section 2.3 and Appendix Section D discuss the construction of the predicted PT spending measure and Section 2 describes the sample definition in further detail. Data: 20% Medicare Carrier claims.

	(1) (2)		(3)	(4)	(5)	(6)			
	Dete	erred	Denied						
	Cond. on	Approach	Cond. on	Attempt	Cond. on Approach				
Age	796	755	.668	.664	.0271	.0534			
	(2.37)	(2.36)	(.613)	(.613)	(2.72)	(2.72)			
Age \times 2006	1.06	1.02	-27.5***	-27.3***	-18.6***	-18.7***			
	(3.84)	(3.85)	(4.54)	(4.31)	(5.06)	(4.91)			
Small provider		-1.18**		.089		752			
		(.55)		(.145)		(.63)			
Small provider \times 2006		905		12.4***		9.87***			
		(.874)		(1.05)		(1.22)			
Outcome mean, 2005	18.6	18.6	1.0	1.0	21.2	21.2			
Outcome mean, 2006	26.5	26.5	22.1	22.1	51.9	51.9			
Week-year FE	Х	Х	Х	Х	Х	Х			
Cluster	Provider	Provider	Provider	Provider	Provider	Provider			
N. Providers	11535	11535	10623	10623	11535	11535			
N. Patients	63238		47919		63238				
N. Observations	72002	72002	103598	103598	72002	72002			

Table A4: Robustness: Regression Results on Screening and Differences byProvider Size using Age, 2005-2006

* p < .10, ** p < .05, *** p < .01. This table presents the coefficients from estimating Equation 3 with log patient age as X_i (columns 1, 3, and 5) and Equation 4 with log patient age as X_i and an indicator for whether a patient goes to an above-median-size provider as D_i (columns 2, 4, and 6). The coefficient β_1 is the estimate for "Age," the coefficient for β_2 is the estimate for "Age $\times 2006$ ", the coefficient for β_3 is the estimate for "Small provider" and β_4 is the estimate for "Small provider $\times 2006$." The regression is clustered at the provider (TIN-state) level. Provider size is measured as a TIN-state's 2006-2008 Medicare patient count and "small provider" denotes a provider with below-median patient count. All specifications restrict to patients over 65. Columns 1-2 and 5-6 restrict to the sample of 2005-2006 patients who ever approach the cap and columns 3-4 restrict to the sample of patients who ever make an attempt to go past the cap, as defined in Section 2.3. Data: 20% Medicare Carrier claims.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	Outcome: Deterred				Outcome	ome: Denied			
	Cond. on Approach			Cor	id. on Atte	mpt	Cond. on Approach		
Predicted Spending	-17***	-17.1***	-17.1***	-1.6***	-1.62***	-1.61***	-20.7***	-20.9***	-20.9***
	(.706)	(.707)	(.722)	(.188)	(.189)	(.194)	(.803)	(.812)	(.834)
Predicted Spending \times 2006	297	0484	.00384	-7.08***	-6.98***	-6.98***	-3.76**	-4.01***	-3.96***
	(1.04)	(1.04)	(1.05)	(1.34)	(1.37)	(1.37)	(1.49)	(1.49)	(1.5)
Non-white	.169			.54**			.892		
	(.614)			(.275)			(.747)		
Non-white \times 2006	-2.72***			8.88***			5.8***		
	(.955)			(1.46)			(1.36)		
Dual		.979*			.437**			1.5^{**}	
		(.578)			(.219)			(.701)	
Dual \times 2006		-2.78***			4.97***			4.37***	
		(.929)			(1.31)			(1.31)	
Low-income subsidy			.323			.401*			.749
			(.569)			(.21)			(.671)
Low-income subsidy \times 2006			-2.45***			4.51***			4.5***
			(.913)			(1.26)			(1.28)
Outcome mean, 2005	18.4	18.4	18.3	1.0	1.0	1.0	21.1	21.1	21.0
Outcome mean, 2006	26.3	26.3	26.3	23.1	23.1	23.1	52.2	52.2	52.2
Week-year FE	Х	Х	Х	Х	Х	Х	Х	Х	Х
Cluster	Provider	Provider	Provider	Provider	Provider	Provider	Provider	Provider	Provider
N. Providers	11911	11911	11887	11064	11064	11041	11911	11911	11887
N. Patients	70518	70518	69788	53560	53560	52986	70518	70518	69788
N. Observations	80532	80532	79782	116360	116360	115228	80532	80532	79782

Table A5: Regression Results on Screening and Differences by Patient Demo-
graphics, 2005-2006

* p < .10, ** p < .05,*** p < .01. This table presents the coefficients from estimating Equation 3 with log predicted PT spending as X_i (columns 1, 3, and 5) and Equation 4 with log predicted PT spending as X_i and an indicator for a patient's demographic characteristic D_i : non-white (columns 1, 4, and 7), dual-eligibility for Medicaid (columns 2, 5, and 8), and Part D Low Income Subsidy status in 2006 (columns 3, 6, and 9). The coefficient β_1 is the estimate for "Predicted spending," the coefficient for β_2 is the estimate for "Predicted spending × 2006", the coefficient for β_3 is the estimate for a demographic indicator and β_4 is the estimate for the demographic indicator, interacted with an indicator for 2006. All specifications include week-year fixed effects for the calendar week(s) that a beneficiary attempts or approaches the cap. The regression is clustered at the provider (TIN-state) level. Columns 1-3 and 7-9 restrict to the sample of 2005-2006 patients who ever approach the cap and columns 4-6 restrict to the sample of patients who ever make an attempt to go past the cap, as defined in Section 2.3. Data: 20% Medicare Carrier claims.

Table A6: Regression Results on Screening and Differences by Patient Demographics with Provider-Year FEs, 2005-2006

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
	Outcome: Deterred			Outcome: Denied								
		Cond. on	Approach			Cond. or	a Attempt			Cond. on	Approach	
Predicted Spending	-13.5*** (.856)	-13.5*** (.856)	-13.6*** (.855)	-13.6*** (.873)	528*** (.147)	528*** (.147)	532*** (.147)	517*** (.148)	-14.5*** (.894)	-14.5^{***} (.894)	-14.6*** (.893)	-14.6*** (.917)
Predicted Spending \times 2006	.492 (1.26)	.488 (1.26)	.521 (1.26)	.545 (1.27)	-4.84*** (.843)	-4.87*** (.842)	-4.85*** (.842)	-4.88*** (.841)	-2.64^{*} (1.4)	-2.64* (1.4)	-2.62* (1.4)	-2.61^{*} (1.42)
Non-white		.874 (.806)		. ,		212 (.189)				.6 (.853)		
Non-white \times 2006		-1.62 (1.2)				1.99^{**}				431 (1.26)		
Dual		()	.867 (.751)			()	.0943 $(.21)$			()	1.15 (.836)	
Dual \times 2006			875 (1.15)				.358				33 (1.29)	
Low-income subsidy			()	.168 $(.741)$			(.0226 $(.199)$			()	.358 (.813)
Low-income subsidy \times 2006				944 (1.14)				.665 (.779)				.43 (1.27)
Outcome mean, 2005	18.0	18.0	18.0	18.0	0.9	0.9	0.9	0.9	20.9	20.9	20.9	20.9
Outcome mean, 2006	24.9	24.9	24.9	24.9	22.4	22.4	22.4	22.4	52.2	52.2	52.2	52.2
Week-year FE	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х
Provider-year FE	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х
Cluster	Provider	Provider	Provider	Provider								
N. Providers	9725	9725	9725	9700	10788	10788	10788	10764	9725	9725	9725	9700
N. Patients	65853	65853	65853	65130	52968	52968	52968	52394	65853	65853	65853	65130
N. Observations	75528	75528	75528	74783	115716	115716	115716	114584	75528	75528	75528	74783

* p < .10, ** p < .05, *** p < .01. This table presents the coefficients from estimating Equation 3 with log predicted PT spending as X_i (columns 1, 5, and 9) and the inclusion of provider-year fixed effects, and Equation 4 with log predicted PT spending as X_i , the inclusion of provider-year fixed effects, and an indicator for a patient's demographic characteristic D_i : non-white (columns 2, 6, and 10), dual-eligibility for Medicaid (columns 3, 7, and 11), and Part D Low Income Subsidy status in 2006 (columns 4, 8, and 12) with provider-year fixed effects. The coefficient β_1 is the estimate for "Predicted spending," the coefficient for β_2 is the estimate for "Predicted spending × 2006", the coefficient for β_3 is the estimate for a demographic indicator and β_4 is the estimate for the demographic indicator, interacted with an indicator for 2006. All specifications include week-year fixed effects for the calendar week(s) that a beneficiary attempts or approaches the cap. The regression is clustered at the provider (TIN-state) level. Columns 1-4 and 9-12 restrict to the sample of 2005-2006 patients who ever approach the cap and columns 5-8 restrict to the sample of patients who ever make an attempt to go past the cap, as defined in Section 2.3. Data: 20% Medicare Carrier claims.

	(1)	(2)						
	Outcome: Denied							
	Cond. on Attempt	Cond. on Approach						
Median predicted spending in 2006 (\$)	1514	1481						
Median week of year	25	25						
Predicted 2006 denial rates								
Small provider	15.8	44.4						
Small - large provider difference	$12.7 \ (80\%)$	9.7~(22%)						
Non-white	20.9	48.1						
Non-white - white difference	8.9~(43%)	5.8 (12%)						
Dual eligible	21.3	48.0						
Dual - non-dual difference	5.0~(23%)	4.4 (9%)						
LIS	21.4	48.1						
LIS - non-LIS difference	4.5~(21%)	4.5 (9%)						

Table A7: Differences by Provider Size, Race, and Income on Denials for MedianPatient

This table presents the implied denial rate in 2006, holding fixed patient predicted spending and the week of approach and attempt, and comparing across race and income groups. The predicted denial rate is calculated using the coefficients from Table A5. Column 1 conditions on patients who make an attempt and column 2 conditions on patients who make an approach, as defined in Section 2.3. Data: 20% Medicare Carrier claims.

B Appendix Figures

Figure B1: Claim-level Denial Rates by Weeks to Cap, Pre- and Post-Reform



(a) Pre-Reform, 2005

(b) Post-Reform, 2006

This figure plots the average denial rate for PT claims depending on how far the patient is from the therapy cap, where the denial rate is the share of claims with at least one denied line. Panel (a) graphs the denial rates relative to a (placebo) 2006 cap in 2005, and panel (b) graphs the denial rates relative to the cap in 2006. A patient's distance to the cap is calculated in terms weeks of care. When "weeks from cap" is negative, it denotes how many additional weeks of care a patient would have to receive to reach the cap. When "weeks from cap" is positive, it denotes how many weeks ago the patient passed the cap. The procedure to define "weeks from cap" is described in detail in Section 2.3. Denial rate is defined as the share of lines in an attempt-week that are denied. Data: 20% Medicare Carrier claims.



Figure B2: Share Receiving Therapy on Each Day of Week

This plot graphs the share of PT visits on each day of the week, given the day of therapy in the previous week. Data: 20% Medicare Carrier claims.





This figure plots the (a) distributions of end-of-year physical therapy spending around the cap in 2004 and 2005 and (b) the difference in the distributions between 2004 to 2005. Distance from cap is calculated in bins of \$50 relative to the 2006 cap and shares are calculated as the share of patients within [-\$800, \$1600] of the cap. Data: 20% Medicare Carrier claims.



Figure B4: Extensive Margin Responses

This figure characterizes potential extensive margin responses to the 2006 therapy cap. Panel (a) plots a kernel density of Medicare physical therapy spending in 2005 and 2006 up to \$3000. Since Medicare pays for 80 percent of allowed charges and patients are responsible for the remaining 20 percent, so the cap appears at $0.8 \times \$1740 = \1392 . Panel (b) plots the average number of beneficiaries per provider, split by whether a provider had a high or low over (placebo) cap share in 2005. Data: 20% Medicare Carrier claims.

Figure B5: Bunching in Dollars, 1999-2000

(a) 1999 and 2000 Spending Distributions (b) Difference between 1999 and 2000



This figure plots the (a) distributions of end-of-year physical therapy spending around the cap in 1999 and 2000 and (b) the difference in the distributions between 1999 and 2000. Distance from cap is calculated in bins of \$50 relative to the 1999 cap (inflation-adjusted to 2005 dollars) and shares are calculated as the share of patients within [-\$700, \$1400] of the cap. Data: 20% Medicare Carrier claims.



Figure B6: Distribution of predicted PT spending, by deterred and denied

This figure plots the distribution predicted 2005 PT spending for "deterred" patients, defined as those who stop care one week before the therapy cap without making an attempt, and "denied" patients, defined as those who stop care the week before the therapy cap after a denied attempt. Panel (a) plots the distribution of predicted spending for each set of patients in 2005, panel (b) plots the distributions in 2006, and panel (c) plots the difference between the 2006 and 2005 distributions, as well as the medians of each distribution (solid line for denied and dotted line for deterred). Section 2.3 describes the construction of the weeks from cap measure as well as the definitions of deterrence and denial in further detail, and Section 4 and Appendix Section D describe the predicted PT spending measure in further detail. Data: 20% Medicare Carrier claims.



Figure B7: Reduced Form Spending and Health Outcomes, Dual Eligibles

This figure plots the coefficient θ_f , which denotes the interaction between an indicator for 2006 and an indicator for month of first PT, from Equation 1. Sample is restricted to beneficiaries with more than \$200 in 12-month PT spending who are dual-eligible for Medicare and Medicaid. Panel (a) plots the coefficients on 12-month PT spending (\$). Panel (b) plots the coefficients on an indicator for pain management procedures, panel (c) plots the coefficients on an indicator for orthopedic surgery, panel (d) plots the coefficients on an indicator for emergency department visit, panel (e) plots the coefficients on an indicator for a hospital stay, panel (f) plots the coefficients on an indicator for a skilled nursing facility stay, and panel (g) plots the coefficients on an indicator for opioid prescriptions. All outcomes other than opioid prescriptions are measured within 12 months of the first PT session; opioid prescriptions are measured 12-24 months after the first PT session. Section E describes how the health outcome measures are defined. Data: 20% Medicare Carrier, Outpatient, MEDPAR, and Part D claims.



Figure B8: Reduced Form Spending and Health Outcomes, Low Income Subsidy

This figure plots the coefficient θ_f , which denotes the interaction between an indicator for 2006 and an indicator for month of first PT, from Equation 1. Sample is restricted to beneficiaries with more than \$200 in 12-month PT spending who receive the Low Income Subsidy in 2006. Panel (a) plots the coefficients on 12-month PT spending (\$). Panel (b) plots the coefficients on an indicator for pain management procedures, panel (c) plots the coefficients on an indicator for orthopedic surgery, panel (d) plots the coefficients on an indicator for emergency department visit, panel (e) plots the coefficients on an indicator for a hospital stay, panel (f) plots the coefficients on an indicator for a skilled nursing facility stay, and panel (g) plots the coefficients on an indicator for opioid prescriptions. All outcomes other than opioid prescriptions are measured within 12 months of the first PT session; opioid prescriptions are measured 12-24 months after the first PT session. Section E describes how the health outcome measures are defined. Data: 20% Medicare Carrier, Outpatient, MEDPAR, and Part D claims.

Figure B9: IV: Effect of PT Spending on Health Outcomes, Low Income Populations



(a) Dual-Eligible Beneficiaries

(b) Low-Income Subsidy Beneficiaries

This figure plots the coefficient β , which denotes the effects of an additional \$100 of PT on an indicator for each outcome, from Equation 2. Panel (a) subsets to dual-eligible beneficiaries and panel (b) subsets to beneficiaries receiving the Part D Low Income Subsidy (LIS). Since LIS is not available until 2006, eligibility for the 2005 sample is based on their 2006 eligibility. All outcomes other than opioid prescriptions are measured within 12 months of the first PT session; opioid prescriptions are measured 12-24 months after the first PT session. Section E describes how the health outcome measures are defined. Sample is restricted to beneficiaries with more than \$200 in 12-month PT spending. Reduced form results are reported in Figures B7 and B8 20% Medicare Carrier, Outpatient, MEDPAR, and Part D claims.

Figure B10: IV: Effect of PT Spending on Health Outcomes, by Decile of Predicted Spending



This figure plots the coefficient β , which denotes the effects of an additional \$100 of PT on an indicator of each outcome, from estimating Equation 2 on different subsets of predicted need deciles. All outcomes other than opioid prescriptions are measured within 12 months of the first PT session; opioid prescriptions are measured 12-24 months after the first PT session. Section D describes how predicted PT spending is constructed, and Section E describes how the health outcome measures are defined. Sample is restricted to beneficiaries with more than \$200 in 12-month PT spending. Data: 20% Medicare Carrier, Outpatient, MEDPAR, and Part D claims.

Figure B11: IV: Effect of PT Spending on Health Outcomes, By Recent Utilization



(a) By Recent Pain Procedure





This figure plots the coefficient β , which denotes the effects of an additional \$100 of PT on an indicator for each outcome, from Equation 2. Results are stratified by whether (a) a patient had a pain procedure in the last 6 months or (b) a patient had an orthopedic surgery in the last 6 months. All outcomes other than opioid prescriptions are measured within 12 months of the first PT session; opioid prescriptions are measured 12-24 months after the first PT session. Section E describes how the health outcome measures are defined. Sample is restricted to beneficiaries with more than \$200 in 12-month PT spending. Reduced form results are reported in Figures B7 and B8 20% Medicare Carrier, Outpatient, MEDPAR, and Part D claims.



(c) Targeting Improvement with Horizontal Inequity



This figure presents illustrations of the screening exercise in Figures 6 and 7. Panel (a) illustrates the case when the introduction of the cap increases the share of patients who stop at the cap (either due to deterrence or denials), but does not improve targeting. The increase in the share stopping at the cap is uniform across all levels of patient need. Panel (b) illustrates the case when the cap increases the share who stop at the cap is larger for low-need patients and smaller for high-need patients. The overall likelihood of receiving care past the cap has become correlated with need once the cap is introduced, indicating that the cap screens on need. Panel (c) illustrates the case when the cap improves targeting and also introduces horizontal inequity on characteristic D. Holding fixed patient need, the increase in share stopping at the cap is larger when D = 1 than when D = 0.

Figure B13: Correlations between patient age and deterrence and denial, 2005-2006



(a) Deterred, conditional on approach

(b) Denied, conditional on attempt



(c) Deterred by provider size, conditional on approach





This figure plots the relationship between attempt outcomes and patient need and provider size. "Deterred" is defined as the share of patients who approach the cap but do not attempt, and "denied" is defined as the share of patients who attempt but never make it past the cap. Sample restricted to patients over the age of 65. Panels (a) and (b) plot the relationship between log age and share deterred and denied in 2005 and 2006. Panels (c) and (d) plot the same relationship, split by provider size. Provider size is defined as the total number of Medicare beneficiaries who receive regular PT by that provider in 2006-2008, and a large firm is defined as being abovemedian. Section 2 describes the sample definition in further detail. Data: 20% Medicare Carrier claims and Master Beneficiary Summary File.


Figure B14: Correlation Between Provider Size and Patient Demographics

This figure plots the coefficient of the relationship between patient characteristics and an indicator variable for going to a small (below-median size) provider, among patients who approach the cap in 2006 (as defined in Section 2.3). Provider size is measured as a TIN-state's 2006-2008 Medicare patient count. Each regression is clustered at the beneficiary-level. Data: 20% Medicare Carrier claims, Master Beneficiary files, and 2006 Individual Income Tax ZIP Code data (SOI Tax Stats, Internal Revenue Service).

Figure B15: Decomposition of Size Advantage on Approvals



(b) Industry-Wide Time Trend (*Week*_t) (c) Provider Fixed Effects (α_i) vs. Size



This figure characterizes differences in approval rates on cap attempts by provider size and experience in 2006-2008 using data and regression estimates from Equation (5). Panel (a) plots the coefficient on the provider's cumulative number of prior attempts on approvals. Panel (b) provides estimates of weekly industry-wide trends. Panel (c) plots the provider fixed effects against provider patient count. Section 2.3 describes the definition of cap attempts in further detail. Data: 20% Medicare Carrier claims.

(a) Learning-by-Doing (κ^e)



Figure B16: IV: Effect of PT Spending on Health Outcomes, By Provider Size

This figure plots the coefficient β , which denotes the effects of an additional \$100 of PT on an indicator for each outcome, from Equation 2. Results are stratified by whether a patient goes to an above-median-sized provider, defined as provider who sees 925 or more Medicare patients in 2006-2008. All outcomes other than opioid prescriptions are measured within 12 months of the first PT session; opioid prescriptions are measured 12-24 months after the first PT session. Section E describes how the health outcome measures are defined. Sample is restricted to beneficiaries with more than \$200 in 12-month PT spending. Reduced form results are reported in Figures B7 and B8 20% Medicare Carrier, Outpatient, MEDPAR, and Part D claims.

Figure B17: Distributions of Spending Around Cap in 2006-2008, Relative to 2005



This figure plots the difference in the distributions of PT spending around the cap between 2005 to (a) 2006 (reproduced from Figure 2), (b) 2007, and (c) 2008. Distance from cap is calculated in bins of \$50 relative to the 2006 cap and shares are calculated as the share of patients within [-\$800, \$1600] of the cap. Data: 20% Medicare Carrier claims.

(a) Excess and Missing Mass in 2006

(b) Excess and Missing Mass in 2007

C Policy Context: 1999 Hard Cap

The first therapy cap regime spanned January-December 1999 and was the result of the Balanced Budget Act of 1997. This legislation introduced two separate \$1500 caps—one for PT/SLP and one for OT.³⁰ The 1999 cap was referred to as a "hard cap" in that there was no exceptions process, and Medicare would not cover any services above the cap. Implementation of the cap was imperfect,³¹ but the cap was highly salient to providers. As a result of aggressive lobbying by the PT industry,³² Congress placed a 2-year moratorium on the cap in 2000, and in subsequent years continued to include provisions to extend the delay of the cap one year at a time.

Figure B5 plots the 1999 and 2000 spending distributions and differences in distributions in the range from \$700 (in 1999 dollars) below the cap to \$1300 above the cap, which is equal to approximately \$800 and \$1600 in 2006 dollars. As depicted in the distributions, some patients do manage to get care reimbursed above the cap in 1999, most likely due to imperfect enforcement of the hard cap. According to a 2003 report on the implementation of the hard cap, CMS could not accurately track cumulative per-patient spending because of "Y2K"-related computing constraints and low provider awareness of the modifier codes introduced to track therapy claims (DynCorp, 2002).

³⁰The three services were intended to each have their own cap, but the combination of PT and SLP into one cap is purportedly the result of a missing "oxford comma" in the text of the legislation (WebPT, 2024)

³¹It was revealed in later reports about the 1999 cap that CMS could not accurately track cumulative per-patient spending because of Y2K-related computing constraints and low provider usage of the modifier codes introduced to track therapy claims.

³²Physical therapy industry representatives reportedly launched "feverish lobbying campaigns ... directed at softening" the payment changes brought on by the Balanced Budget Act (Goldstein, 1999). The American Physical Therapy Association recruited thousands of physical therapists to stage protests on Capitol Hill and organized phone-a-thons to call on Congressional representatives. (Luthra, 2018).

D Machine Learning Methodology and Results

Patient-level 12 Month PT Spending Prediction In order to construct a proxy for patient need, we construct a patient-level measure of predicted PT spending based on patient characteristics and utilization in the 6 months and the year *prior* to starting PT. We then apply this prediction to all 2005-2006 patients in the health analysis sample in Section 3.2 and patients who approach and attempt the cap in Section 4. The model is trained on patients who approach or attempt the cap in 2004 and 2005, prior to the implementation of the therapy cap. We use gradient-boosted decision trees from the LightGBM package. The predictors are: age, race, sex, utilization and spending in the previous calendar year available in the MBSF Cost and Utilization file (in-office spending, Part B drug, outpatient procedure, inpatient, testing, imaging, hospice, evaluation and management, durable medical equipment, dialysis, and other), chronic conditions at the end of the previous calendar year, PT and OT spending in the previous calendar year, inpatient and SNF stays within the last 6 months (spending, number of visits, Diagnosis Related Group in most recent visit, length of stay of last visit, and days since last visit), in-office spending in the last 5 months.



Figure D1: Predicted vs. Actual 12-month Spending, 2004-2005

This figure plots the predicted 12-month physical therapy spending against the actual 12-month spending in 2004-2005 from the model described in Section 2.3 and Appendix Section D. Data: 20% Medicare Carrier claims and Master Beneficiary Summary File.

Figure D1 plots the relationship between actual 12-month PT spending and predicted spending based on the model; there is a monotonic relationship between predicted and actual spending and the R^2 of the prediction is over 10 percent. The predictors with the highest feature importance are (in order of importance): Part B physician office spending in the previous year, Part B drug spending in the previous year, total in-office spending in the last 6 months, patient age, spending on tests in the previous year, durable medical equipment spending in the previous year, total outpatient spending in the last 6 months, and the number of imaging events in the previous year.

Claim-level Denial Rate Prediction Table 1 showed that including documentation substantially increases the explanatory power of a regression model in explaining denials. However, this regression model is unable to capture the high-dimensional information contained on claims, such as procedure codes and diagnosis codes. Thus, it may not be able to capture complex billing behavior like upcoding which could be used to increase approval

rates. To explore whether strategic billing behavior could be driving approvals over time, we train a machine learning model that uses line- and claim-level information from claims associated with the attempt as well as the weeks of care leading up to the attempt. We first use data from 2006 to train a machine learning model that predicts the likelihood of claim-level denial rate for claims aiming to exceed the cap using information from the claim. We again use gradient-boosted decision trees from the LightGBM package. The predictors include patient age, race, sex, ICD-9 diagnosis codes on the claim, the HCPCS procedure code on the line item, the number of units on that line item, modifier codes on the line item (including the KX modifier), number of units, modifier codes on the patient's last 5 visits, as well as the prior year and prior 6 month utilization and spending used in the 12 month PT spending prediction described above. Intuitively, this model is a probabilistic approximation to the decisions made by the Medicare contractor in deciding denials. We then apply the prediction to claims associated with attempts in 2005 to 2008, setting the KX modifier to be on or off for all lines.





This figure plots claim-level predicted denial rates from the denial rate prediction model described in Section 5. Panel (a) plots the predicted rates for claims associated with attempts in 2006, split by whether the claim was actually denied or not. Panels (b) predicted rates, split by whether the claim had documentation or not. The prediction model is trained on 2006 claims associated with cap attempts and explanatory variables in the prediction model are discussed in Appendix Section D. Section 2.3 describes the definition of cap attempts in further detail. Data: 20% Medicare Carrier claims and Master Beneficiary Summary File.

Figure D2a shows that our model predicts denials out of sample well. For denied claims, the median predicted probability of denial is 79%. We label 90 percent of denied claims as having more than 50% probability of denial. Likewise, for approved claims, the median predicted denial probability is just 16%. We label 78 percent of these approved claims as having less than 50% probability of denial.

While the model incorporates many more variables than the regression in Table 1, again we find that one factor has outsize predictive influence: documentation. Figure D2b shows that even though the machine learning model was trained using a host of patient- and claimlevel characteristics, the documentation indicator alone explains a substantial fraction of the prediction. Indeed, this is because documentation explains a large amount of variation in the data: in 2006, just 37% of claims using documentation were denied, while 75% of claims without documentation were denied.

Finally, we investigate whether the change in approval rates over time could reflect changes in strategic billing behavior. In order to capture billing behavior that is independent of documentation usage, we use our model to create a predicted approval rate under the assumption that *every attempt uses documentation*. The time series of this modified predicted approval rate in Figure D3 shows that changes in billing behavior besides documentation are unable to account for most of the increase in the approval rate over time – if all claims were coded as having documentation, the approval rates would have largely persisted at their pre-cap levels.

Figure D3: Time Series Variation in Documentation and Predicted Approval Rate on Cap Attempts



This figure plots the share of attempts with no denials ("Share approved") and predicted approval rate based on claim characteristics, assuming documentation ("Predicted approval rate with documentation"). Attempts are defined as weeks in which the patient attempts to go past the cap, where an attempt is as defined in Section 2.3: a week in which a patient's cumulative paid amount is below the cap at the beginning of the week, and their cumulative paid or billed amount is above the cap at the end of the week. The predicted approval rate is derived from the line-level denial rate prediction described in Section D on claims associated with attempts, where the KX modifier indicator for documentation use is turned on for all claims. 20% Medicare Carrier claims.

E Construction of Patient Health Measures

Injury diagnoses were identified as claims within 90 days of the PT start date in MEDPAR (inpatient hospital stay), Outpatient, and Carrier files an ICD-9 code starting with 800-899. Pain management procedures were identified in the Outpatient and Carrier files based on HCPCS codes. Orthopedic surgeries were identified in the Outpatient files using based on HCPCS codes, respectively. Crosswalks to HCPCS codes for pain procedures and orthopedic surgeries were created in consultation with a clinical expert (available upon request). Emergency department visits were identified as inpatient stays in MEDPAR with positive ED charges or Outpatient claims with a Revenue Center Code between 450 and 459 or equal to

0981. Hospital stays were identified in MEDPAR as claims with Short Stay/Long Stay/SNF Provider Indicator Code "S" or "L," and skilled nursing stays were identified as claims with Short Stay/Long Stay/SNF Provider Indicator Code "N." Opioid prescriptions were identified in the Part D file as prescriptions with Product Service ID/National Drug Service Code via a crosswalk to Anatomical Therapeutic Chemical codes for opioids, in consultation with a clinical expert (available upon request).

F Patient Health Effects: Difference-in-Difference Strategy

Sample and Identification Strategy In addition to the identification strategy in Section 3.2 which leverages within-year differences in spending depending on when a patient starts PT, we deploy an alternative empirical strategy to assess the patient health effects of the cap. This empirical strategy is in the spirit of the estimator used in Diamond and Persson (2016), which proposes a method to estimate the causal effects of bunching by comparing outcomes for individuals inside and outside of a "manipulation region" around a discontinuity. In our case, we compare average health outcomes among patients whose end-of-year spending could have plausibly have been reduced by the cap, who serve as the treated group, and patients whose spending is too low to have been affected by the cap, who serve as the control group. Figure G1 illustrates how we define the two groups. The treated group includes any patients who are over the cap or within 5 weeks of the cap, while the control group includes all patients over 5 weeks under the cap.³³ The underlying assumption is that the control group does not include any "bunchers" whose spending was at risk of being affected by the cap.

 $^{^{33}}$ Because we convert a patient's end-of-year spending into "weeks from cap" using the maximum of their 5-week rolling average spending *or* the sample average weekly spending, this left-censors the "weeks from cap" measure at -8. The -8 week bin can be interpreted as patients who appear to end the year 8 or more weeks below the cap. We also restrict to patients who are at most 48 weeks to the right of the cap to exclude patients who are implausibly far from the cap. We also restrict to patients who end their PT outside of the first and last 4 weeks of the year.

Figure G1: Patient Health Outcomes DD Identification Strategy



This figure illustrates the treatment (gray) and control (white) group assignment for the differencein-difference patient health identification strategy described in Section F.

If the reductions in spending from the cap led to worsened patient health, then once the cap is in place, we would expect the treated group's *average* health to fall relative to the control group. To see this, consider a patient who *would* have ended up above the cap, but instead bunches under the cap once it is in place. This patient contributes to the treated group's average both before and after the reform. Thus, if reducing spending harmed this patient, we would expect the average health for the treated group should fall. In contrast, the control group is comprised of patients who would remain far below the cap regardless of whether the cap is in place. We do not have to know *which* patients below the cap are there because of the cap, but rather just that some share of them would have counterfactually been above the cap.

After constructing the treatment and control groups, we use a difference-in-difference strategy to compare how their health outcomes evolve before and after the 2006 cap. Interpreting the estimates from this difference-in-difference as causal requires making two assumptions. The first is the standard parallel trends assumption: the health trends for each group would be parallel in the absence of the reform. We will verify this by looking for evidence of pre-trends in an event study. The second assumption is that the composition of patients in each group is unaffected by the cap. In other words, none of the "bunchers" reduced their spending so much that they ended up over 5 weeks away from the cap. Figure 3 shows that there is no statistically significant difference in the 2005 and 2006 distributions past 4 weeks from the cap. Additionally, the lack of an extensive margin response in Figure B4 suggests that the cap did not affect the composition of patients who seek PT. Furthermore, in our main specification we directly control for several observable patient characteristics.

Results We present the results in the form of a yearly event study. The specification for patient *i* receiving care in $Year_i$ that ends in calendar week $LastWeek_i \in [1, 52]$ is:

$$Y_i = \beta_0 + \beta_1 Treated_i + \sum_{\tau=2004}^{2008} \beta_{2\tau} Treated_i \times 1(Year_i = \tau) + LastWeek_i + \varepsilon_i.$$
(7)

The pooled specification is:

$$Y_i = \beta_0 + \beta_1 Treated_i + \beta_2 Post_i + LastWeek_i + \varepsilon_i, \tag{8}$$

where $Post_i$ is an indicator for years after 2006. The estimates are clustered at the patient level and the omitted year is 2005. Figure G2 plots the β_{2t} estimates from Equation 7 and Figure G3 plots the coefficients the pooled specification in Equation 8. Figure G2 panel (a) confirms a "first stage"—the treated group saw a sizeable reduction in PT spending relative to the control group. Turning to health outcomes, consistent with the null results found in Section 3.2, we find no evidence of increases in (b) the usage pain management procedures, (c) orthopedic surgeries, (d) hospital stays, or (e) skilled nursing home stays.³⁴

 $^{^{34}}$ Given that the Medicare Part D program for prescription medication began in 2006, we cannot use this difference-in-difference strategy to study opioid prescriptions.

Figure G2: DD Design: Spending and Health Outcomes within 12 Months of First PT



This figure plots the coefficients from the health outcomes regression in Equation 7. Panel (a) plots the coefficients on log 12-month PT spending. All outcomes are all measured within 12 months of first PT session. Panel (b) plots the coefficients on an indicator for pain management procedures, panel (c) plots the coefficients on an indicator for orthopedic surgery, panel (d) plots the coefficients on an indicator for an emergency department visit, panel (e) plots the coefficients on an indicator for a skilled nursing facility stay. The left-side y-axis denotes the effect in terms of percentage points and the right-side y-axis denotes the effect in terms of percent of the control group comprises patients who end the year within 5 weeks away from the cap or above the cap, and the control group compromises patients who end the year over 5 weeks below the cap. Section 2.3 provides additional detail on how "weeks from the cap" is defined. Data: 20% Medicare Carrier claims and Master Beneficiary Summary Files. 85

Figure G3: DD Design: Spending and Health Outcomes within 12 Months of First PT



This figure plots the coefficients from the health outcomes regression in Equation 8. All outcomes are all measured within 12 months of first PT session in a given year. The treated group comprises patients who end the year within 5 weeks away from the cap or above the cap, and the control group compromises patients who end the year over 5 weeks below the cap. Section 2.3 provides additional detail on how "weeks from the cap" is defined. Data: 20% Medicare Carrier claims and Master Beneficiary Summary Files.